

Pairing interacting protein sequences using masked language modeling

Umberto Lupo^{1,2,*†}, Damiano Sgarbossa^{1,2,*}, Anne-Florence Bitbol^{1,2,†}

1 Institute of Bioengineering, School of Life Sciences, École Polytechnique
Fédérale de Lausanne (EPFL), CH-1015 Lausanne, Switzerland

2 SIB Swiss Institute of Bioinformatics, CH-1015 Lausanne, Switzerland

* These authors contributed equally to this work.

† Emails: umberto.lupo@epfl.ch, anne-florence.bitbol@epfl.ch

Abstract

Predicting which proteins interact together from amino-acid sequences is an important task. We develop a method to pair interacting protein sequences which leverages the power of protein language models trained on multiple sequence alignments, such as MSA Transformer and the EvoFormer module of AlphaFold. We formulate the problem of pairing interacting partners among the paralogs of two protein families in a differentiable way. We introduce a method called DiffPALM that solves it by exploiting the ability of MSA Transformer to fill in masked amino acids in multiple sequence alignments using the surrounding context. MSA Transformer encodes coevolution between functionally or structurally coupled amino acids within protein chains. It also captures inter-chain coevolution, despite being trained on single-chain data. Relying on MSA Transformer without fine-tuning, DiffPALM outperforms existing coevolution-based pairing methods on difficult benchmarks of shallow multiple sequence alignments extracted from ubiquitous prokaryotic protein datasets. It also outperforms an alternative method based on a state-of-the-art protein language model trained on single sequences. Paired alignments of interacting protein sequences are a crucial ingredient of supervised deep learning methods to predict the three-dimensional structure of protein complexes. Starting from sequences paired by DiffPALM substantially improves the structure prediction of some eukaryotic protein complexes by AlphaFold-Multimer. It also achieves competitive performance with using orthology-based pairing.

Significance statement

Deep learning has brought major advances to the analysis of biological sequences. Self-supervised models, based on approaches from natural language processing and trained on large ensembles of protein sequences, efficiently learn statistical dependence in this data. This includes coevolution patterns between structurally or functionally coupled amino acids, which allows them to capture structural contacts. We propose a method to pair interacting protein sequences which leverages the power of a protein language model trained on multiple sequence alignments. Our method performs well for small datasets that are challenging for existing methods. It can improve structure prediction of protein complexes by supervised methods, which remains more challenging than that of single-chain proteins.

Introduction

Interacting proteins play key roles in cells, ensuring the specificity of signaling pathways and forming multi-protein complexes that act e.g. as molecular motors or receptors. Predicting protein-protein interactions and the structure of protein complexes are important questions in computational biology and biophysics. Indeed, high-throughput experiments capable of resolving protein-protein interactions remain challenging [1], even for model organisms, and experimental determination of protein complex structure is demanding.

A major advance in protein structure prediction was achieved by AlphaFold [2] and other deep learning approaches [3–5]. Extensions to protein complexes have been proposed [6–9], including AlphaFold-Multimer (AFM) [7], but their performance is heterogeneous and less impressive than for monomers [10]. Importantly, the first step of AlphaFold is to build multiple-sequence alignments (MSAs) of homologs of the query protein sequence. The results of the CASP15 structure prediction contest demonstrated that MSA quality is crucial to further improving AlphaFold performance [11, 12]. For protein complexes involving several different chains (heteromers), paired MSAs, whose rows include actually interacting chains, can provide coevolutionary signal between interacting partners that is informative about inter-chain contacts [13–16]. However, constructing paired MSAs poses the challenge of properly pairing sequences. Accordingly, the quality of pairings strongly impacts the accuracy of heteromer structure prediction [9, 17, 18]. Pairing interaction partners is difficult because many protein families contain several paralogous proteins encoded within the same genome. This problem is known as paralog matching. In prokaryotes, genomic proximity can often be used to solve it, since most interaction partners are encoded in close genomic locations [19, 20]. However, this is not the case in eukaryotes. Large-scale coevolution studies of protein complexes [21–23] and deep learning approaches [6–9, 24] have paired sequences by using genomic proximity when possible [8, 21, 24], and/or by pairing together the closest, or equally ranked, hits to the query sequences, i.e. relying on approximate orthology [6–9, 17, 22–25].

Aside from genomic proximity and orthology, phylogeny-based methods have addressed the paralog matching problem [26–35], exploiting similarities between evolutionary histories of interacting proteins [36–40]. Other methods, based on coevolution [13, 41–45], rely on correlations in amino-acid usage between interacting proteins [14, 15, 46, 47]. These correlations arise from the need to maintain physico-chemical complementarity among amino acids in contact, and from shared evolutionary history [48, 49]. Phylogeny and coevolution can be explicitly combined, improving performance [50]. However, coevolution-based approaches are data-thirsty and need large and diverse MSAs to perform well. This limits their applicability, especially to eukaryotic complex structure prediction. Nevertheless, the core idea of finding pairings that maximise coevolutionary signal holds promise for paralog matching and complex structure prediction.

We develop a new coevolution-based method for paralog matching which leverages recent neural protein language models taking MSAs as inputs [2, 51]. These models are one of the ingredients of the success of AlphaFold [2]. We focus on MSA Transformer [51], a protein language model which was trained on MSAs using the masked language modeling (MLM) objective in a self-supervised way. We introduce DiffPALM, a differentiable method for predicting paralog matchings using MLM. We show that it outperforms existing coevolution methods by a large margin on difficult benchmarks of shallow MSAs extracted from ubiquitous prokaryotic protein datasets. DiffPALM performance further quickly improves when known interacting pairs are provided as examples. Next, we apply DiffPALM to the hard problem of paralog matching for eukaryotic protein complexes. For this, we take sequences paired by DiffPALM as input to AFM. Among the complexes we tested, using DiffPALM substantially improves structure prediction by AFM in some cases, and, when using orthologs as known interacting pairs, does

not yield any significant deterioration. It also achieves competitive performance with using orthology-based pairing.

Results

Leveraging MSA-based protein language models for paralog matching

MSA-based protein language models, which include MSA Transformer [51] and the EvoFormer module of AlphaFold [2], are trained to correctly fill in masked amino acids in MSAs with the MLM loss (see “[Supplementary material, MSA Transformer and masked language modeling for MSAs](#)” and [51, 52] for details). To this end, they use the rest of the MSA as context, which allows them to capture coevolution. Indeed, MSA Transformer achieves state-of-the-art performance at unsupervised structural contact prediction [51], captures pairwise phylogenetic relationships between sequences [52], and can be used to generate new sequences from given protein families [53]. While MSA Transformer was only trained on MSAs corresponding to single chains, inter-chain coevolutionary signal has strong similarities with intra-chain signal [13–15]. As illustrated by [Fig. S1](#), MSA Transformer is able to detect inter-chain contacts from a properly paired MSA [54]. [Fig. S1](#) further shows that it cannot do so from a wrongly paired MSA. Moreover, we find that the MLM loss (used for the pre-training of MSA Transformer) decreases as the fraction of correctly matched sequences increases, see [Fig. S2](#). These results demonstrate that MSA Transformer captures inter-chain coevolutionary signal.

In this context, we ask the following question: Can we exploit MSA-based protein language models to address the paralog matching problem? Let us focus on the case where two MSAs have to be paired, which is the relevant one for heterodimers. Paralog matching amounts to pairing these two MSAs, each corresponding to one of two interacting protein families, so that correct interaction partners are placed on the same row of the paired MSA. Throughout, we will assume that interactions are one-to-one, excluding cross-talk, which is valid for proteins that interact specifically [55]. Thus, within each species, assuming that there is the same number of sequences from both families, we aim to find the correct one-to-one matching that associates one protein from the first family to one protein from the second family. We also cover the case where the two protein families have different numbers of paralogs within the same species, see “[Methods](#)”. Motivated by our finding that the MLM loss is lower for correctly paired MSAs than for incorrectly paired ones, we address the paralog matching problem by looking for pairings that minimise an MLM loss. A challenge is that the number of possible such one-to-one matchings scales factorially with the number of sequences in the species, making it difficult to find the permutation that minimises the loss by a brute-force search. We address this challenge by formulating a differentiable optimization problem that can be solved using gradient methods, to yield configurations minimizing our MLM loss, see “[Methods](#)”. We call our method **DiffPALM**, short for **D**ifferentiable **P**airing using **A**lignment-based **L**anguage **M**odels.

DiffPALM outperforms other coevolution methods on small MSAs

We start out by considering a well-controlled benchmark dataset composed of ubiquitous prokaryotic proteins from two interacting families, namely histidine kinases (HKs) and response regulators (RRs) [56, 57], see “[Supplementary material, Datasets](#)”. These proteins interact together within prokaryotic two-component signaling systems, important pathways that enable bacteria to sense and respond to environment signals [55]. They possess multiple paralogs (on the order of ten per genome, with substantial variability), and have strong specificity for their cognate partners. Because most cognate HK-RR pairs are encoded in the same operon, many

interaction partners are known from genome proximity, which enables us to assess performance. In addition, earlier coevolution methods for paralog matching were tested on this dataset, allowing rigorous comparison [14, 47, 50]. Here, we focus on datasets comprising about 50 cognate HK-RR pairs. Indeed, this small data regime is problematic for existing coevolution methods, which require considerably deeper alignments to achieve good performance [14, 15, 47, 50]. Furthermore, this regime is highly relevant for eukaryotic complexes, because their homologs have relatively low sequence diversity, as shown by the effective depth of their MSAs in [Table S1](#). While prokaryotic proteins such as HKs and RRs feature high diversity, focusing on small datasets allows us to address the relevant regime of low diversity in this well-controlled benchmark case. We hypothesize that MSA Transformer’s extensive pre-training can help to capture coevolution even in these difficult cases. To assess this, we first test two variants of our DiffPALM method (see “[Methods](#)”) on 40 MSAs from the HK-RR dataset comprising about 50 HK-RR pairs each (see “[Supplementary material, Datasets](#)”). We first address the *de novo* pairing prediction task, starting from no known HK-RR pair, and then we study the impact of starting from known pairs.

[Fig. 1](#) shows that DiffPALM performs better than the chance expectation, obtained for random within-species matching. Moreover, it outperforms other coevolution-based methods, namely DCA-IPA [14], MI-IPA [47], which rely respectively on Potts models and on mutual information, and GA-IPA [50], which combines these coevolution measures with sequence similarity, a proxy for phylogeny. Importantly, these results are obtained without giving any paired sequences as input to the algorithm. The performance of DiffPALM is particularly good for pairs with high confidence score (see “[Result and confidence](#)”), as shown by the “precision-10” curve, which focuses on top 10% predicted pairs, when ranked by predicted confidence (see [Fig. 1](#)). We also propose a method based on a protein language model trained on a large ensemble of single sequences, ESM-2 (650M) [5], see “[Pairing based on a single-sequence language model](#)”. DiffPALM also outperforms this method, even though the latter is faster (no need for backpropagation) and is formulated as a linear matching problem, which is solved exactly. This confirms that the coevolution information contained in the MSA plays a key role in the performance of DiffPALM, which is based on MSA Transformer. A key strength of MSA Transformer and thus of DiffPALM is that they leverage the power of large language models while starting from MSAs, and thus allow direct access to the coevolutionary signal. [Fig. 1](#) shows that both variants of DiffPALM, namely Multi-Run Aggregation (MRA) and Iterative Pairing Algorithm (IPA), outperform all baselines, and that precision of MRA increases with the number of runs used (see [Table S2](#) for details). In MRA, we aggregate pairs predicted via independent optimization runs, while in IPA, we gradually add pairs with high confidence as positive examples (see “[Methods](#)”). DiffPALM-IPA thus exploits the good results obtained for high-confidence pairs, evidenced by the high precision-10 scores in [Fig. 1](#). Note that the distribution of precision-10 scores over the different MSAs we considered is skewed, especially after many MRA runs, see [Fig. S3](#). For many MSAs, almost perfect scores are reached, while performance is bad for a few others. MSAs with smaller average number of sequence per species tend to yield larger precision, as the pairing task is then easier.

So far, we addressed *de novo* pairing prediction, where no known HK-RR pair is given as input. Can DiffPALM precision increase by exploiting “positive examples” of known interacting partners? This is an important question, since experiments on model species may for instance give some positive examples (see “[The paralog matching problem](#)”). To address it, we included different numbers of positive examples, by using the corresponding non-masked interacting pairs as context (see “[Construction of an appropriate MLM loss](#)”). The left panel of [Fig. 2](#) shows that the performance of DiffPALM significantly increases with the number of positive examples used, reaching almost perfect performance for the highest-confidence pairs (precision-10).

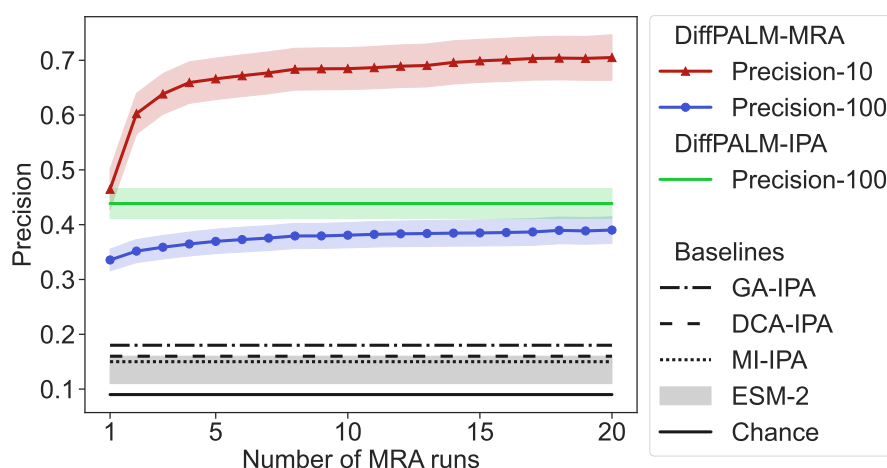


Figure 1: Performance of DiffPALM on small HK-RR MSAs. The performance of two variants of DiffPALM (MRA and IPA, see “[Improving precision: MRA and IPA](#)”) is shown versus the number of runs used for the MRA variant, for 40 MSAs comprising about 50 HK-RR pairs. The chance expectation, and the performance of various other methods, are reported as baselines. Three existing coevolution-based methods are considered: DCA-IPA [14], MI-IPA [47], and GA-IPA [50]. We also consider a pairing method based on the scores given by the ESM-2 (650M) single-sequence protein language model [5], see “[Pairing based on a single-sequence language model](#)”. With all methods, a full one-to-one within-species pairing is produced, and performance is measured by precision (also called positive predictive value or PPV), namely, the fraction of correct pairs among predicted pairs. The default score is “precision-100”, where this fraction is computed over all predicted pairs (100% of them). For DiffPALM-MRA, we also report “precision-10”, which is calculated over the top 10% predicted pairs, when ranked by predicted confidence within each MSA (see “[Methods](#)”). For DiffPALM, we plot the mean performance on all MSAs (color shading), and the standard error range (shaded region). For our ESM-2 based method, we consider 10 different values of masking probability p from 0.1 to 1.0, and we report the range of precisions obtained (gray shading). For other baselines, we report the mean performance on all MSAs.

Another important parameter is MSA depth. Indeed, deeper MSAs yield better performance for traditional coevolution methods [14, 47, 50]. The middle panel of Fig. 2 shows that DiffPALM performance also increases with MSA depth, while already reaching good performance for shallow MSAs. Note that MSA Transformer’s large memory footprint makes it difficult to consider very deep input MSAs (see “[Supplementary material, Datasets](#)”).

While we focused on HK-RR pairing so far, DiffPALM is a general method. To assess how it extends to other cases, we consider another pair of ubiquitous prokaryotic proteins, namely homologs of the *E. coli* proteins MALG-MALK, which are involved in ABC transporter complexes. These proteins form permanent complexes, while HK-RR interact transiently to transmit signal. The right panel of Fig. 2 shows results obtained on 40 MSAs comprising about 50 MALG-MALK pairs, without positive examples. We observe that DiffPALM outperforms the chance expectation by a large margin. It also significantly outperforms existing coevolution methods [14, 47, 50], as well as our method based on ESM-2 (650M), see Table S2. Note that all approaches yield better performance for MALG-MALK than for HK-RR, as the number of MALG-MALK pairs per species is smaller than that of HK-RR pairs. Finally, while Fig. 2 reports the final MRA performance, Fig. S4 shows that both performance scores increase with the number of MRA runs.

HK-RR and MALG-MALK are interesting benchmarks because they have multiple paralogs per species and because ground-truth pairings are known from genome proximity. However,

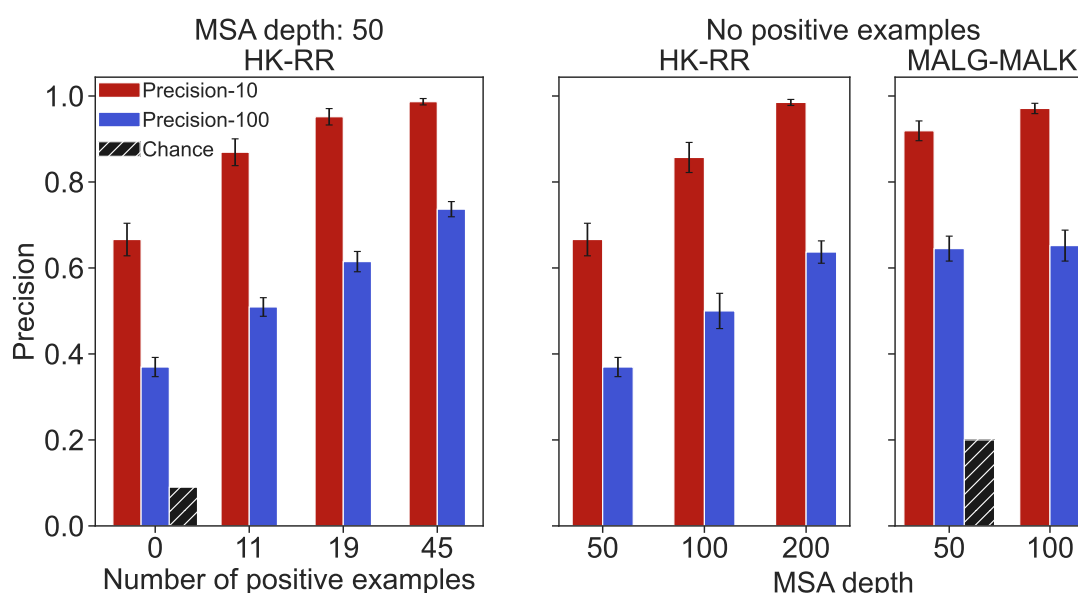


Figure 2: Impact of positive examples, MSA depth, and extension to another pair of protein families. We report the performance of DiffPALM with 5 MRA runs (measured as precision-100 and precision-10, see Fig. 1), for various numbers of positive examples, on the same HK-RR MSAs as in Fig. 1 (left panel). We also report the performance of DiffPALM (using no positive examples) versus MSA depth for both HK-RR and MALG-MALK pairs (middle and right panel). In all cases, we show the mean value over different MSAs and its standard error, and we plot the chance expectation for reference. Note that MSA depth can vary by $\pm 10\%$ around the reported value because complete species are used (see “Supplementary material, Datasets”).

both domains involved in these pairs are sometimes present within a single chain. In fact, gene fusions often exist for interacting proteins [58], and are used as a clue to predict interactions [59]. While hybrid HKs comprising both HK- and RR-specific domains are known to have lower specificity and weaker coevolutionary signal than other pairs [45], the presence of such proteins in the training set of MSA Transformer could be an ingredient of the success of DiffPALM. However, such proteins are also included in the training set of ESM-2, which we found to perform far less well than MSA Transformer for pairing. To further investigate if their presence in the training set of MSA Transformer is key to DiffPALM performance, we considered another prokaryotic pair, NUOA-NUOJ (NADH-quinone oxidoreductase). This system has very few fusions: only one protein with both domains is referenced in the InterPro database [60]. We ran DiffPALM for 4 MSAs of 50 pairs of NUOA-NUOJ and obtained an average precision-100 score of 0.86 (standard error 0.06), to be compared to a chance expectation of 0.5. Thus, DiffPALM performs well even when gene fusions are very rare.

Using DiffPALM for eukaryotic complex structure prediction by AFM

An important and more challenging application of DiffPALM is predicting interacting partners among the paralogs of two families in eukaryotic species. Indeed, eukaryotes often have many paralogs per species [61] but eukaryotic-specific protein families generally have fewer total homologs and smaller diversity than prokaryotes. Moreover, most interacting proteins are not encoded in close proximity in eukaryotic genomes. Paired MSAs are a key ingredient of protein complex structure prediction by AFM [7, 9]. When presented with query sequences, the default AFM pipeline [7] retrieves homologs of each of the chains. Within each species, homologs of different chains are ranked according to Hamming distance to the corresponding

query sequence. Then, equal-rank sequences are paired. Can DiffPALM improve complex structure prediction by AFM? To address this question, we consider 15 complexes, listed in [Table S1](#), whose structures are not included in the training set of the AFM release we used unless specified otherwise (v2, see “[Supplementary material, General points on AFM](#)”), and for which the default AFM complex prediction was previously reported to perform poorly [7, 18] (see “[Supplementary material, Eukaryotic complexes](#)”).

[Fig. 3](#) (top panels) shows that DiffPALM can improve complex structure prediction by AFM (see [Fig. S5](#) for details). This suggests that it is able to produce better paired MSAs than those from the default AFM pipeline. In particular, substantial improvements are obtained for the complexes with PDB identifiers 6L5K and 6FYH, see [Figs. S6](#) and [S7](#) for structural visualizations. For 6FYH, using DiffPALM also yields much higher average AFM confidence scores than the default AFM pipeline ([Fig. 3](#), bottom panels). For 6L5K, both pairing methods yield very high confidence scores, but none of the structures predicted with the default AFM pipeline agree with the experimental structure. We found no cases in which using DiffPALM improved AFM confidence scores but deteriorated structure predictions.

In most cases, the quality of structures predicted using DiffPALM pairing is comparable to that obtained using the pairing method adopted e.g. by ColabFold [8], where only the orthologs of the two query sequences, identified as their best hits, are paired in each species (resulting in at most one pair per species) [8, 9, 22–25], see [Fig. 3](#). Note however that, for 6PNQ, the ortholog-only pairing method is outperformed both by DiffPALM and by the default AFM pairing. Indeed, the raw and effective MSA depths are smaller for this structure than e.g. for 6L5K and 6FYH (see [Table S1](#)). Thus, further reducing diversity by restricting to paired orthologs may be negatively impacting structure prediction in this case. Given the good results obtained overall with orthology-based pairings, we tried using them as positive examples for DiffPALM. Given the very good precision obtained by DiffPALM for high-confidence HK-RR pairs (see above), we also tried restricting to high-confidence pairs. For most structures, we obtained no significant improvement over the standard DiffPALM using these variants (see [Fig. S8](#)). However, for 6WCW, we could generate several higher-quality structures, particularly when using orthologs as positive examples. Overall, when compared with the default AFM pairing pipeline, DiffPALM with orthologs as positive examples never significantly deteriorates prediction quality for our structures.

Although DiffPALM achieves similar performance on these structure prediction tasks as using orthology, it predicts some pairs that are quite different from orthology-based pairs. Indeed, [Fig. S9](#) shows that the fraction of pairs identically matched by DiffPALM and by orthology is often smaller than 0.5. [Fig. S10](#) further shows that, for the sequences that are paired differently by DiffPALM and by orthology, the Hamming distances between the two predicted partners is often above 0.5. Nevertheless, most of the pairs that are predicted both by DiffPALM and by using orthology have high DiffPALM confidence (see [Fig. S9](#)), confirming the importance of these pairs.

How does DiffPALM compare to traditional coevolution methods for these eukaryotic complexes? To address this, we paired chains using MI-IPA [47] for the 6 complexes with deepest pairable MSAs (see [Table S1](#)), and used these paired MSAs as input for AFM. [Fig. S11](#) shows that MI-IPA also yields some improvement over default AFM pairing for the structures improved by DiffPALM. However, DiffPALM yields stronger improvement than MI-IPA in these cases. This confirms the efficiency of DiffPALM at using coevolutionary signal to pair the sequences. Note that we focused on deep MSAs because of the depth requirements of traditional coevolution methods, and that DiffPALM performed better despite deep MSAs being split for DiffPALM due to memory limitations (see “[Supplementary material, Datasets](#)”), which reduces the coevolutionary information available to DiffPALM.

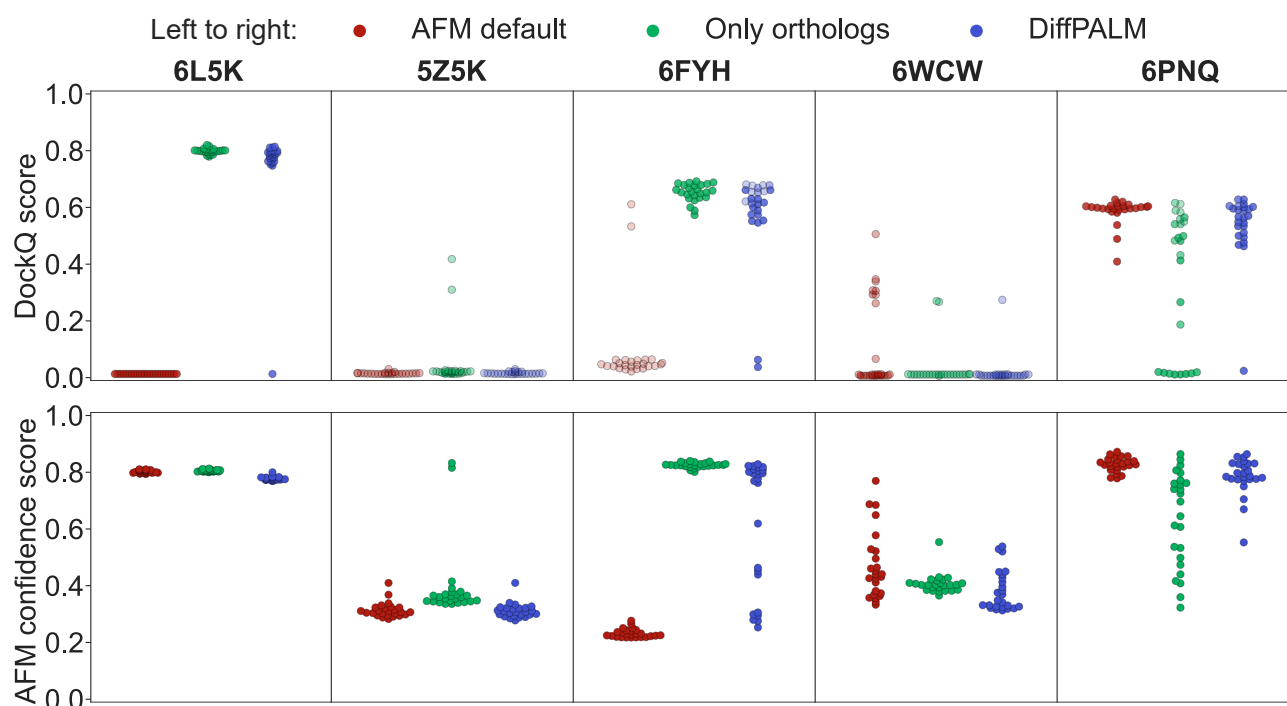


Figure 3: Performance of AFM using different pairing methods. We use AFM to predict the structure of protein complexes starting from differently paired MSAs, each of them constructed from the same initial unpaired MSAs. Three pairing methods are considered: the default one of AFM, only pairing orthologs to the two query sequences, and a single run of DiffPALM (equivalent to one MRA run). We used a single run for computational time reasons. Performance is evaluated using DockQ scores (top panels), a widely used measure of quality for protein-protein docking [62], and the AFM confidence scores (bottom panels), see “[Supplementary material, General points on AFM](#)”. The latter are also used as transparency levels in the top panels, where more transparent markers denote predicted structures with low AFM confidence. For each query complex, AFM is run five times. Each run yields 25 predictions which are ranked by AFM confidence score. The top five predicted structures are selected from each run, giving 25 predicted structures in total for each complex. Out of the 15 complexes listed in [Table S1](#), we restrict to those where any two of these three pairing methods yield a significant difference (> 0.1) in average DockQ scores for at least one set of predictions coming from different runs but with the same within-run rank according to AFM confidence. Panels are ordered by increasing mean DockQ score for the AFM default method.

Finally, [Fig. S12](#) compares the performance of DiffPALM to the AFM default and the ortholog-only pairing methods using the latest AFM release (v3, see “[Supplementary material, Eukaryotic complexes](#)”). Since all the structures considered here are included in the training set of this release, we expect the standard AFM pipeline to perform well. However, we observe that using DiffPALM yields similar or better performance than using the AFM default or ortholog-only pairing methods. In particular, the structure prediction of 6L5K remains substantially improved by using DiffPALM rather than AFM default pairing. Moreover, for 6THL, which was poorly predicted by all methods with the older AFM release (DockQ < 0.1), DiffPALM yields a substantial improvement compared to both other pairing methods with the latest AFM release.

Discussion

We developed DiffPALM, a method for pairing interacting protein sequences that builds on MSA Transformer [51], a protein language model trained on MSAs. MSA Transformer efficiently captures coevolution between amino acids, thanks to its training to fill in masked amino acids using the surrounding MSA context [51–53]. It also captures inter-chain coevolutionary signal, despite being trained on single-chain MSAs [54]. We leveraged this ability in DiffPALM by using a masked language modeling loss as a coevolution score and looking for the pairing that minimizes it. We formulated the pairing problem in a differentiable way, allowing us to use gradient methods. On shallow MSAs extracted from controlled prokaryotic benchmark datasets, DiffPALM outperforms existing coevolution-based methods and a method based on a state-of-the-art language model trained on single sequences. Its performance quickly increases when adding examples of known interacting sequences. It also increases with MSA depth, as for traditional coevolution methods. Paired MSAs of interacting partners are a key ingredient to complex structure prediction by AFM. We found that using DiffPALM can improve the performance of AFM, and achieves competitive performance with orthology-based pairing.

Recent work [18] also used MSA Transformer for paralog matching, in a method called ESMPair. It relies on column attention matrices and compares them across the MSAs of interacting partners. This makes it quite different from DiffPALM, which relies on coevolutionary information via the MLM loss. ESMPair may be more closely related to phylogeny- or orthology-based pairing methods, since column attention encodes phylogenetic relationships [52]. 13 out of the 15 eukaryotic protein complexes we considered were also studied in [18], but no substantial improvement (and often a degradation of performance) was reported for those when using ESMPair instead of the default AFM pairing, except for 7BQU. By contrast, DiffPALM yields strong improvements for 6L5K and 6FYH, and no significant performance degradation when using orthologs as positive examples. Explicitly combining coevolution and phylogeny using MSA Transformer is a promising direction to further improve partner pairing. Indeed, such an approach has already improved traditional coevolution methods [50]. Other ways of improving MSA usage by AFM have also been proposed [63] and could be combined with advances in pairing. Besides improving MSA construction [64] and the extraction of MSA information, other promising approaches include exploiting structural alignments [65], using massive sampling and dropout [66], and combining AFM with more traditional docking methods [67, 68], which has allowed e.g. to improve structure prediction of 6A6I [67].

While DiffPALM performance increases with MSA depth, MSA Transformer’s memory and time requirements scale quadratically with MSA depth and with MSA length, limiting the size of input MSAs. Using e.g. FlashAttention [69], or using a simpler MSA-based protein language model in place of MSA Transformer, might help alleviate some of these limitations.

DiffPALM illustrates the power of neural protein language models trained on MSAs, and their ability to capture the rich structure of biological sequence data. The fact that these models encode inter-chain coevolution, while they are trained on single-chain data, suggests that they are able to generalize. We used MSA Transformer in a zero-shot setting, without fine-tuning it to the task of interaction partner prediction. Such fine-tuning could yield further performance gains [70].

Like MSA Transformer, AlphaFold’s EvoFormer [2] also processes MSA input. Furthermore, a portion of the total loss used during its training as part of a supervised structure prediction method consisted of an MLM loss. Thus, one could develop a paralog matching method based on EvoFormer. However, the substantially larger pre-training dataset of MSA Transformer should make DiffPALM more broadly applicable. Furthermore, the monomer version of EvoFormer performs less well than MSA Transformer for mutational effect prediction, despite being an excellent contact predictor [71]. Since mutational effect prediction probes the

quality of the entire probability vectors output at masked positions, this hints that DiffPALM might not perform as well if we were to use the monomer version of EvoFormer instead of MSA Transformer. Nevertheless, an interesting perspective would be to integrate paralog matching into an end-to-end retraining of AlphaFold Multimer.

The fact that DiffPALM outperforms existing coevolution methods [14, 47, 50] for shallow MSAs is reminiscent of the impressive performance of MSA Transformer at predicting structural contacts from shallow MSAs [51]. While traditional coevolution methods either compute local coevolution scores for two columns of an MSA [47] or build a global model for an MSA [14, 15], MSA Transformer was trained on large ensembles of MSAs and shares parameters across them. This presumably allows it to transfer knowledge between MSAs, and to bypass the usual needs for deep MSAs of traditional coevolution methods [14, 15, 46, 47, 50], or of MSA-specific transformer models [72]. This constitutes major progress for the use of coevolutionary signal.

After the transformative progress brought by deep learning to protein structure prediction [2–5], predicting protein complex structure and ligand binding sites is fast advancing with AFM and related methods, but also with other deep learning models based on structural representations [73–76]. Combining the latter [77, 78] and, more generally, structural information [79] with the power of sequence-based language models is starting to bring even further progress.

Methods

The paralog matching problem

Goal and notations. Paralog matching amounts to pairing a pair of MSAs, each one corresponding to one of the two protein families considered. We assume that interactions are one-to-one. Let $\mathcal{M}^{(A)}$ and $\mathcal{M}^{(B)}$ be the (single-chain) MSAs of two interacting protein families A and B, and let K denote the number of species represented in both MSAs and comprising more than one unmatched sequence in at least one MSA. Species represented in only one MSA are discarded since no within-species matching is possible for them. Species with only one unmatched sequence in each MSA are not considered further since pairing is trivial. There may also be N_{pos} known interacting pairs: they are treated separately, as positive examples (see below). Here we focus on the unmatched sequences. For $k = 1, \dots, K$, let $N_k^{(A)}$ and $N_k^{(B)}$ denote the number of unmatched sequences belonging to species k in $\mathcal{M}^{(A)}$ and $\mathcal{M}^{(B)}$ (respectively).

Dealing with asymmetric cases. The two protein families considered may have different numbers of paralogs within the same species. Assume, without loss of generality, that $N_k^{(A)} < N_k^{(B)}$ for a given k . To solve the matching problem with one-to-one interactions, we would like to pick, for each of the $N_k^{(A)}$ sequences in $\mathcal{M}^{(A)}$, a single and exclusive interaction partner out of the $N_k^{(B)}$ available sequences in $\mathcal{M}^{(B)}$. The remaining sequences of the species in $\mathcal{M}^{(B)}$ are left unpaired. In practice, we achieve this by augmenting the original set of species- k sequences from $\mathcal{M}^{(A)}$ with $N_k^{(B)} - N_k^{(A)}$ “padding sequences” made entirely of gap symbols. By doing so (and analogously when $N_k^{(A)} > N_k^{(B)}$), the thus-augmented interacting MSAs have the same number $N_k := \max(N_k^{(A)}, N_k^{(B)})$ of sequences from each species k . In practice, this method is used for the AFM complex structure prediction, while the curated benchmark prokaryotic MSAs do not have asymmetries (see “[Supplementary material, Datasets](#)”).

Formalization. The paralog matching problem corresponds to finding, within each species k , a mapping that associates one sequence of $\mathcal{M}^{(A)}$ to one sequence of $\mathcal{M}^{(B)}$ (and reciprocally). Thus, within each species k , one-to-one matchings can be encoded as permutation matrices of size $N_k \times N_k$. A brute-force search through all possible within-species one-to-one matchings

would scale factorially with the size N_k of each species, making it prohibitive. Note that the Iterative Pairing Algorithm (IPA) [14, 47] is an approximate method to solve this problem when optimizing coevolution scores. Here, we introduce another one, which allows to leverage the power of deep learning.

Exploiting known interacting partners. Our use of a language model allows for contextual conditioning, a common technique in natural language processing. Indeed, if any correctly paired sequences are already known, they can be included as part of the joint MSA input to MSA Transformer. In this case, we exclude their pairing from the optimization process – in particular, by not masking any of their amino acids, see below. We call these known paired sequences “positive examples”. In Fig. 2, we randomly sampled species and included all their pairs as positive examples, until we reached the desired depth $N_{\text{pos}} \pm 10\%$. For eukaryotic complex structure prediction, we treated the query sequence pair as a positive example.

DiffPALM: Paralog matching based on MLM

Here, we explain our paralog matching method based on MLM, which we call DiffPALM. Background information on MSA Transformer and its MLM loss is collected in “Supplementary material, MSA Transformer and masked language modeling for MSAs”. DiffPALM exploits our differentiable framework for optimizing matchings, see “Supplementary material, A differentiable formulation of paralog matching”. The key steps are summarized in Fig. 4.

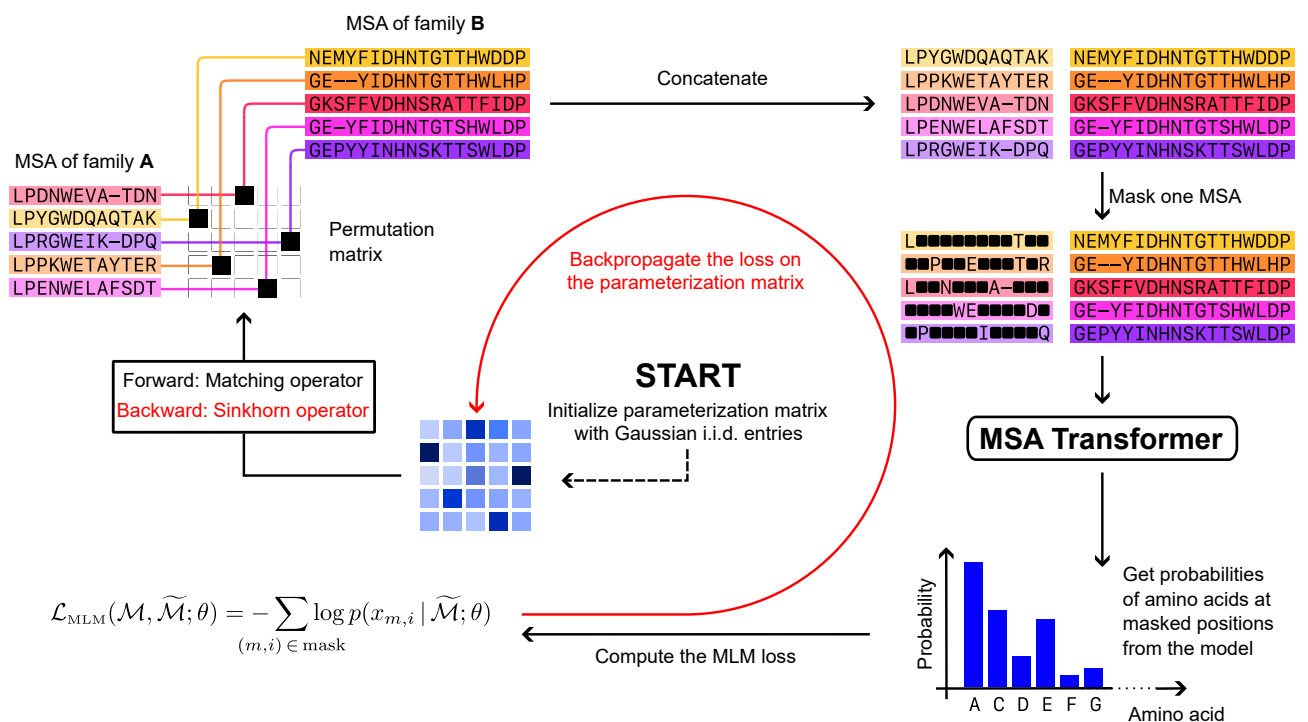


Figure 4: **Schematic of the DiffPALM method.** First, the parameterization matrices X_k are initialized, and then the following steps are repeated until the loss converges: (1) Compute the permutation matrix $M(X_k)$ and use it to shuffle $\mathcal{M}^{(A)}$ relative to $\mathcal{M}^{(B)}$. Then pair the two MSAs. (2) Randomly mask some tokens of one of the two sides of the paired MSA and compute the MLM loss Eq. (S1). (3) Backpropagate the loss and update the parameterization matrices X_k , using the Sinkhorn operator \hat{S} for the backward step instead of the matching operator M (see “Supplementary material, A differentiable formulation of paralog matching”).

Construction of an appropriate MLM loss. Using the tools just described, we consider two interacting MSAs (possibly augmented with padding sequences), still denoted by $\mathcal{M}^{(A)}$ and $\mathcal{M}^{(B)}$. Given species indexed by $k = 1, \dots, K$, we initialize a set $\{X_k\}_{k=1, \dots, K}$ of square matrices of size $N_k \times N_k$ (the case $K = 1$ corresponds to X in “[Supplementary material, A differentiable formulation of paralog matching](#)”). We call these “parameterization matrices”. By applying to them the matching operator M [see [Eq. \(S2\)](#)], we obtain the permutation matrices $\{M(X_k)\}_{k=1, \dots, K}$, encoding matchings within each species in the paired MSA. Using gradient methods, we optimize the parameterization matrices so that the corresponding permutation matrices yield a paired MSA with low MLM loss. More precisely, paired MSAs are represented as concatenated MSAs with interacting partners placed on the same row,¹ and our MLM loss for this optimization is computed as follows:

1. Perform a shuffle of $\mathcal{M}^{(A)}$ relative to $\mathcal{M}^{(B)}$ using the permutation matrix $M(X_k)$ in each species k (plus an optional noise term, see below), to obtain a paired MSA \mathcal{M} ;
2. Generate a mask for \mathcal{M} (excluding any positive example tokens from the masking);
3. Compute MSA Transformer’s MLM loss for that mask, see [Eq. \(S1\)](#).

Importantly, we only mask tokens from one of the two MSAs, chosen uniformly at random within that MSA with a high masking probability $p \geq 0.7$.² Our rationale for using large masking probabilities is that it forces the model to predict masked residues in one of the two MSAs by using information coming mostly from the other MSA – see [Fig. S2](#). We stress that, if padding sequences consisting entirely of gaps are present (see above), we mask these symbols with the same probability as those coming from ordinary sequences. Of the two MSAs to pair, we mask the one with shorter length if no padding sequences exist (i.e. here for our prokaryotic benchmark datasets). Else, if lengths are comparable but one MSA contains considerably more padding sequences than the other, we preferentially mask that MSA. Otherwise, we randomly choose which of the two MSAs to mask.

We fine-tuned all the hyperparameters involved in our algorithm using two joint MSAs of depth ~ 50 , constructed by selecting all the sequences of randomly sampled species from the HK-RR dataset (see “[Supplementary material, Datasets](#)”).

Noise and regularization. Following [80], after updating (or initializing) each X_k , we add to it a noise term given by a matrix of standard i.i.d. Gumbel noise multiplied by a scale factor. The addition of noise ensures that the X_k do not get stuck at degenerate values for the right-hand side of [Eq. \(S2\)](#), and more generally may encourage the algorithm to explore larger regions in the space of permutations. As scale factor for this noise we choose 0.1 times the sample standard deviation of the current entries of X_k , times a global factor tied to the optimizer scheduler (see next paragraph). Finally, since the matching operator is scale-invariant, we can regularize the matrices X_k to have small Frobenius norm. We find this to be beneficial and implement it through weight decay, set to be $w = 0.1$.

Optimization. We backpropagate the MLM loss on the parameterization matrices X_k . After each gradient step, we mean-center the parameterization matrices. We use the AdaDelta optimizer [81] with an initial learning rate $\gamma = 9$ and a “reduce on loss plateau” learning rate scheduler which decreases the learning rate by a factor of 0.8 if the loss has not decreased for more than 20 gradient steps after the learning rate was last set. The learning rate scheduler also provides the global scale factor which, together with the standard deviation of the entries of X_k , dynamically determines the magnitude of the Gumbel noise (see previous paragraph).

¹We employ no special tokens to demarcate the boundary between one sequence and its partner.

²In contrast, uniformly random masking with $p = 0.15$ was used during MSA Transformer’s pre-training [51].

Exploring the loss landscape through multiple initializations. We observe that the initial choice of the parameterization set $\{X_k\}_{k=1,\dots,K}$ strongly impact results. Slightly different initial conditions for X_k lead to very different final permutation matrices. Furthermore, we observe fast decrease in the loss when the X_k are initialized to be exactly zero (our use of Gumbel noise means that we break ties randomly when computing the permutation matrices $M(X_k)$; if noise is not used, similar results can be achieved by initializing X_k with entries very close to zero). Thus, we can cheaply probe different paths in the loss landscape by performing several short runs using zero-initialized parameterization matrices X_k . In practice, we use 20 different such short runs each consisting of 20 gradient steps. Then, we average all the final parameterizations together to warm-start a longer run made up of 400 gradient steps.

Result and confidence. We observe that, even though the loss generally converges to a minimum average value during our optimization runs, there are often several distinct hard permutations associated to the smallest loss values. This may indicate a flattening of the loss landscape relative to the inherent fluctuations in the MLM loss, and/or the existence of multiple local minima. To extract a single matching per species from one of our runs (or indeed from several runs, see “[Improving precision: MRA and IPA](#)”), we average the hard permutation matrices associated to the q lowest losses, and evaluate the matching operator [Eq. (S2)] on the resulting averages. We find final precision-100 figures to be quite robust to the choice of q . On the other hand, for individual (warm-started) runs as described in “[Exploring the loss landscape through multiple initializations](#)”, precision-10 benefits from setting q to its maximum possible value of 400.

Furthermore, we propose using each entry in the averaged permutation matrices as an indicator of the model’s confidence in the matching of the corresponding pair of sequences. Indeed, pairs that are present in most low-loss configurations are presumably essential for the optimization process and are captured most of the times, pushing their confidence value close to 1. Conversely, non-interacting pairs are in most of the cases associated to higher losses and therefore appear sporadically, obtaining confidences close to zero. Accordingly, we refer to the averaged hard permutations used to extract a single matching per species as “confidence matrices”, and to the final in-species matchings as “consensus permutations”.

Improving precision: MRA and IPA. We propose two methods for improving precision further. In the first method, which we call Multi-Run Aggregation (MRA), we perform N_{runs} independent optimization runs for each interacting MSA. Then, we collect together the hard permutations independently obtained from each run, and aggregate the $q = 400$ lowest-loss permutations from this larger collection to obtain more reliable confidence matrices and hard permutations.

The second method is an iterative procedure analogous to the Iterative Pairing Algorithm (IPA) of Refs. [14, 47], and named after it. The idea is to gradually add pairs with high confidence as positive examples. Assuming a paired MSA containing a single species for notational simplicity, the n -th iteration (starting at $n = 1$) involves the following steps:

1. Perform an optimization run and extract from it a confidence matrix $C^{(n)}$ as described in “[Result and confidence](#)”, using the currently available positive examples;
2. Compute the moving average $\tilde{C}^{(n)} = \text{mean}(C^{(n)}, \tilde{C}^{(n-1)}, \dots, \tilde{C}^{(1)})$ (where $\tilde{C}^{(1)} \equiv C^{(1)}$);
3. Define candidate matchings via the consensus permutation $M^{(n)} = M(\tilde{C}^{(n)})$;
4. Repeat Steps 1-3 a maximum of 3 times, until the average MLM loss estimated using $M^{(n)}$, and 200 random masks, is lower or statistically insignificantly higher³ than what could have been obtained using $M^{(n-1)}$ and the same positive examples as in Step 1;

³Based on a two-sample T-test: we say “statistically insignificantly” when $p \geq 0.95$, and “statistically significantly” when $p < 0.05$.

5. If Step 4 fails, set $\tilde{C}^{(n)} = \tilde{C}^{(n-1)}$ and $M^{(n)} = M(\tilde{C}^{(n-1)})$ (but removing rows and columns corresponding to the positive examples added at iteration $n - 1$);
6. Check that the average MLM loss estimated using $M^{(n)}$ and 200 random masks, but only regarding as positive examples those available at the beginning of iteration $n - 1$, is not statistically significantly higher³ than the average MLM loss estimated using $M^{(n-1)}$ and those same positive examples;
7. If Step 6 fails, terminate the IPA. Otherwise, pick the top 5 pairs according to $\tilde{C}^{(n)}$, promote them to positive examples in all subsequent iterations, and remove them from the portion of the paired MSA to be optimized.

If several species are present, they are optimized together (see “[Construction of an appropriate MLM loss](#)”) and confidence values from all species are used to select the top 5 pairs.

Pairing based on a single-sequence language model

To assess whether a single-sequence model is able to solve the paralog matching problem, we consider the 650M-parameter version of the model ESM-2 [5]. We score candidate paired sequences using the MLM loss in Eq. (S1). In contrast with MSA Transformer, the input of the model is not paired MSAs but single paired sequences. Therefore, it is sufficient to individually score each possible pair within each species, without needing to consider all permutations. Denoting by N_k the number of sequences from each family in species k , the number of possible pairs is N_k^2 while the number of permutations is $N_k!$. This complexity reduction allows us to evaluate the scores of all possible pairs. This removes the need of backpropagating the loss on the permutation matrix. Accordingly, this method is much faster, since we only need to use the model in evaluation mode, without gradient backpropagation.

For each candidate paired sequence, we evaluate the average of the MLM losses computed over multiple random masks (with masking probability p). Once the average MLM losses are computed for all the N_k^2 pairs, we compute the optimal one-to-one matching by using standard algorithms for linear assignment problems [82] on the $N_k \times N_k$ matrix containing all the losses.

Assessing the impact of pairing on AFM structure prediction

Pairing methods employed in AFM and ColabFold. When presented with a set of query chains, AFM retrieves homologs of each of the chains by running JackHMMER [83] on UniProt, and further homology searches on other databases [7]. UniProt hits are partitioned into species⁴ and ranked within each species by decreasing Hamming distance to the relevant query sequence. A paired MSA is obtained by matching equal-rank hits. Sequences left unpaired are discarded. In addition, AFM produces “block MSAs” constructed by “pairing” hits from the remaining databases with padding sequences of gaps. The input for AFM comprises the paired MSA and the block MSAs.

While sharing the same architecture and weights as AFM, ColabFold retrieves homologs using MMseqs2 [84] on ColabFoldDB [8]. In each species, hits are sorted by increasing E-value, and the best hits are paired [8, 9, 22–25]. Thus, contrary to the default AFM pipeline, the paired MSA in ColabFold contains at most one sequence pair per species for a heterodimer. Because the databases and homology search methods used by ColabFold differ from those used by AFM, a direct comparison does not allow one to isolate the effect of their different pairing schemes. Therefore, we employed the ColabFold pairing method starting from the sequences that are paired in the default AFM pipeline.

⁴While the official AFM implementation in <https://github.com/deepmind/alphafold> uses UniProt “entry names” to define species, when possible we instead use NCBI TaxIDs (via UniProt mappings to NCBI tax IDs, retrieved on 17 Dec 2022, corresponding to UniProt Knowledgebase release 2022_04), which are more accurate.

Pairing using DiffPALM. To assess the impact of DiffPALM on complex structure prediction by AFM, we started from the sequences that are paired in the default AFM pipeline. We left out species with large unbalances between the number of sequences in the two families considered. Specifically, if the ratio of the larger to the smaller of these two numbers exceeds an ad-hoc “maximum size ratio” MSR (see Table S1), if there is only one sequence in both families, or if there are more than 50 sequences in at least one family, then we do not attempt pairing via DiffPALM, and revert to default AFM pairing. When the full MSA to be paired with DiffPALM is too deep and/or long, optimizing it as a whole is not possible due to GPU memory limitations. Instead, we partition it into several small enough sub-MSAs, which we optimize independently. Note that this may reduce performance, since pairing quality increases with MSA depth. We always use the query sequences as positive examples. For all these optimization runs, we use a masking probability $p = 0.7$.

Acknowledgments

The authors thank Sergey Ovchinnikov for valuable feedback on a preliminary version of this manuscript, and Alex Hawkins-Hooker for interesting conversations. This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No. 851173, to A.-F. B.).

Data availability statement

A Python implementation of DiffPALM is freely available in our GitHub repository: <https://github.com/Bitbol-Lab/DiffPALM>.

References

- [1] S. V. Rajagopala, P. Sikorski, A. Kumar, R. Mosca, J. Vlasblom, R. Arnold, J. Franca-Koh, S. B. Pakala, S. Phanse, A. Ceol, R. Hauser, G. Siszler, S. Wuchty, A. Emili, M. Babu, P. Aloy, R. Pieper, and P. Uetz, “The binary protein-protein interaction landscape of *Escherichia coli*,” *Nat Biotechnol*, vol. 32, no. 3, pp. 285–290, 2014.
- [2] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko, A. Bridgland, C. Meyer, S. A. A. Kohl, A. J. Ballard, A. Cowie, B. Romera-Paredes, S. Nikolov, R. Jain, J. Adler, T. Back, S. Petersen, D. Reiman, E. Clancy, M. Zielinski, M. Steinegger, M. Pacholska, T. Berghammer, S. Bodenstein, D. Silver, O. Vinyals, A. W. Senior, K. Kavukcuoglu, P. Kohli, and D. Hassabis, “Highly accurate protein structure prediction with AlphaFold,” *Nature*, vol. 596, no. 7873, pp. 583–589, 2021.
- [3] M. Baek, F. DiMaio, I. Anishchenko, J. Dauparas, S. Ovchinnikov, G. R. Lee, J. Wang, Q. Cong, L. N. Kinch, R. D. Schaeffer, C. Millán, H. Park, C. Adams, C. R. Glassman, A. DeGiovanni, J. H. Pereira, A. V. Rodrigues, A. A. van Dijk, A. C. Ebrecht, D. J. Opperman, T. Sagmeister, C. Buhlheller, T. Pavkov-Keller, M. K. Rathinaswamy, U. Dalwadi, C. K. Yip, J. E. Burke, K. C. Garcia, N. V. Grishin, P. D. Adams, R. J. Read, and D. Baker, “Accurate prediction of protein structures and interactions using a three-track neural network,” *Science*, vol. 373, no. 6557, pp. 871–876, 2021.
- [4] R. Chowdhury, N. Bouatta, S. Biswas, C. Floristean, A. Kharkar, K. Roy, C. Rochereau, G. Ahdritz, J. Zhang, G. M. Church *et al.*, “Single-sequence protein structure prediction

- using a language model and deep learning,” *Nat Biotechnol*, vol. 40, no. 11, pp. 1617–1623, 2022.
- [5] Z. Lin, H. Akin, R. Rao, B. Hie, Z. Zhu, W. Lu, N. Smetanin, R. Verkuil, O. Kabeli, Y. Shmueli, A. dos Santos Costa, M. Fazel-Zarandi, T. Sercu, S. Candido, and A. Rives, “Evolutionary-scale prediction of atomic-level protein structure with a language model,” *Science*, vol. 379, no. 6637, pp. 1123–1130, 2023.
- [6] I. Humphreys, J. Pei, M. Baek, A. Krishnakumar, I. Anishchenko, S. Ovchinnikov, J. Zhang, T. J. Ness, S. Banjade, S. R. Bagde, V. G. Stancheva, X. H. Li, K. Liu, Z. Zheng, D. J. Barrero, U. Roy, J. Kuper, I. S. Fernández, B. Szakal, D. Branzei, J. Rizo, C. Kisker, E. C. Greene, S. Biggins, S. Keeney, E. A. Miller, J. C. Fromme, T. L. Hendrickson, Q. Cong, and D. Baker, “Computed structures of core eukaryotic protein complexes,” *Science*, vol. 374, no. 6573, 2021.
- [7] R. Evans, M. O’Neill, A. Pritzel, N. Antropova, A. Senior, T. Green, A. Židek, R. Bates, S. Blackwell, J. Yim, O. Ronneberger, S. Bodenstern, M. Zielinski, A. Bridgland, A. Potapenko, A. Cowie, K. Tunyasuvunakool, R. Jain, E. Clancy, P. Kohli, J. Jumper, and D. Hassabis, “Protein complex prediction with AlphaFold-Multimer,” *bioRxiv*, 2021.
- [8] M. Mirdita, K. Schütze, Y. Moriwaki, L. Heo, S. Ovchinnikov, and M. Steinegger, “ColabFold: making protein folding accessible to all,” *Nat Methods*, vol. 19, no. 6, pp. 679–682, 2022.
- [9] P. Bryant, G. Pozzati, and A. Elofsson, “Improved prediction of protein-protein interactions using AlphaFold2,” *Nat Commun*, vol. 13, no. 1, p. 1265, 2022.
- [10] H. Schweke, T. Levin, M. Pacesa, C. A. Goverde, P. Kumar, Y. Duhoo, L. J. Dornfeld, B. Dubreuil, S. Georgeon, S. Ovchinnikov, D. N. Woolfson, B. E. Correia, S. Dey, and E. D. Levy, “An atlas of protein homo-oligomerization across domains of life,” *bioRxiv*, 2023.
- [11] A. Elofsson, “Progress at protein structure prediction, as seen in CASP15,” *Current Opinion in Structural Biology*, vol. 80, p. 102594, 2023.
- [12] L. T. Alexander, J. Durairaj, A. Kryshchak, L. A. Abriata, Y. Bayo, G. Bhabha, C. Breyton, S. G. Caulton, J. Chen, S. Degroux, D. C. Ekiert, B. S. Erlandsen, P. L. Fredolino, D. Gilzer, C. Greening, J. M. Grimes, R. Grinter, M. Gurusaran, M. D. Hartmann, C. J. Hitchman, J. R. Keown, A. Kropp, P. Kursula, A. L. Lovering, B. Lemaitre, A. Lia, S. Liu, M. Logotheti, S. Lu, S. sson, M. D. Miller, G. Minasov, H. H. Niemann, F. Opazo, G. N. Phillips, O. R. Davies, S. Rommelaere, M. Rosas-Lemus, P. Roversi, K. Satchell, N. Smith, M. A. Wilson, K. L. Wu, X. Xia, H. Xiao, W. Zhang, Z. H. Zhou, K. Fidelis, M. Topf, J. Moult, and T. Schwede, “Protein target highlights in CASP15: Analysis of models by structure providers,” *Proteins*, pp. 1–29, 2023.
- [13] M. Weigt, R. A. White, H. Szurmant, J. A. Hoch, and T. Hwa, “Identification of direct residue contacts in protein-protein interaction by message passing,” *Proc. Natl. Acad. Sci. U.S.A.*, vol. 106, no. 1, pp. 67–72, 2009.
- [14] A.-F. Bitbol, R. S. Dwyer, L. J. Colwell, and N. S. Wingreen, “Inferring interaction partners from protein sequences,” *Proc. Natl. Acad. Sci. U.S.A.*, vol. 113, no. 43, pp. 12 180–12 185, 2016.

- [15] T. Gueudre, C. Baldassi, M. Zamparo, M. Weigt, and A. Pagnani, “Simultaneous identification of specifically interacting paralogs and interprotein contacts by direct coupling analysis,” *Proc. Natl. Acad. Sci. U.S.A.*, vol. 113, no. 43, pp. 12 186–12 191, 2016.
- [16] H. Szurmant and M. Weigt, “Inter-residue, inter-protein and inter-family coevolution: bridging the scales,” *Current Opinion in Structural Biology*, vol. 50, pp. 26–32, 2018.
- [17] Y. Si and C. Yan, “Protein complex structure prediction powered by multiple sequence alignments of interologs from multiple taxonomic ranks and AlphaFold2,” *Briefings in Bioinformatics*, vol. 23, no. 4, p. bbac208, 2022.
- [18] B. Chen, Z. Xie, J. Qiu, Z. Ye, J. Xu, and J. Tang, “Improved the heterodimer protein complex prediction with protein language models,” *Briefings in Bioinformatics*, vol. 24, no. 4, p. bbad221, 2023.
- [19] S. Orchard, M. Ammari, B. Aranda, L. Breuza, L. Briganti, F. Broackes-Carter, N. H. Campbell, G. Chavali, C. Chen, N. del Toro, M. Duesbury, M. Dumousseau, E. Galeota, U. Hinz, M. Iannuccelli, S. Jagannathan, R. Jimenez, J. Khadake, A. Lagreid, L. Licata, R. C. Lovering, B. Meldal, A. N. Melidoni, M. Milagros, D. Peluso, L. Perfetto, P. Porras, A. Raghunath, S. Ricard-Blum, B. Roechert, A. Stutz, M. Tognolli, K. van Roey, G. Cesareni, and H. Hermjakob, “The MIntAct project—IntAct as a common curation platform for 11 molecular interaction databases,” *Nucleic Acids Research*, vol. 42, no. D1, pp. D358–D363, 2013.
- [20] J. M. Peters, A. Colavin, H. Shi, T. L. Czarny, M. H. Larson, S. Wong, J. S. Hawkins, C. H. Lu, B.-M. Koo, E. Marta, A. L. Shiver, E. H. Whitehead, J. S. Weissman, E. D. Brown, L. S. Qi, K. C. Huang, and C. A. Gross, “A comprehensive, CRISPR-based functional analysis of essential genes in bacteria,” *Cell*, vol. 165, no. 6, pp. 1493–1506, 2016.
- [21] S. Ovchinnikov, H. Kamisetty, and D. Baker, “Robust and accurate prediction of residue-residue interactions across protein interfaces using evolutionary information,” *eLife*, vol. 3, p. e02030, 2014.
- [22] Q. Cong, I. Anishchenko, S. Ovchinnikov, and D. Baker, “Protein interaction networks revealed by proteome coevolution,” *Science*, vol. 365, no. 6449, pp. 185–189, 2019.
- [23] A. G. Green, H. Elhabashy, K. P. Brock, R. Maddamsetti, O. Kohlbacher, and D. S. Marks, “Large-scale discovery of protein interactions at residue resolution using co-evolution calculated from genomic sequences,” *Nat Commun*, vol. 12, no. 1, p. 1396, 2021.
- [24] H. Zeng, S. Wang, T. Zhou, F. Zhao, X. Li, Q. Wu, and J. Xu, “ComplexContact: a web server for inter-protein contact prediction using deep learning,” *Nucleic Acids Research*, vol. 46, no. W1, pp. W432–W437, 2018.
- [25] G. Pozzati, W. Zhu, C. Bassot, J. Lamb, P. Kundrotas, and A. Elofsson, “Limits and potential of combined folding and docking,” *Bioinformatics*, vol. 38, no. 4, pp. 954–961, 2021.
- [26] C.-S. Goh and F. E. Cohen, “Co-evolutionary analysis reveals insights into protein–protein interactions,” *Journal of Molecular Biology*, vol. 324, no. 1, pp. 177–192, 2002.
- [27] A. K. Ramani and E. M. Marcotte, “Exploiting the co-evolution of interacting proteins to discover interaction specificity,” *Journal of Molecular Biology*, vol. 327, no. 1, pp. 273–284, 2003.

- [28] J. Gertz, G. Elfond, A. Shustrova, M. Weisinger, M. Pellegrini, S. Cokus, and B. Rothschild, “[Inferring protein interactions from phylogenetic distance matrices](#),” *Bioinformatics*, vol. 19, no. 16, pp. 2039–2045, 2003.
- [29] J. M. Izarzugaza, D. Juan, C. Pons, J. A. Ranea, A. Valencia, and F. Pazos, “[TSEMA: interactive prediction of protein pairings between interacting families](#),” *Nucleic Acids Research*, vol. 34, no. Web Server issue, pp. W315–319, 2006.
- [30] E. R. Tillier, L. Biro, G. Li, and D. Tillo, “[Codep: maximizing co-evolutionary interdependencies to discover interacting proteins](#),” *Proteins*, vol. 63, no. 4, pp. 822–831, 2006.
- [31] J. M. Izarzugaza, D. Juan, C. Pons, F. Pazos, and A. Valencia, “[Enhancing the prediction of protein pairings between interacting families using orthology information](#),” *BMC Bioinformatics*, vol. 9, p. 35, 2008.
- [32] E. R. Tillier and R. L. Charlebois, “[The human protein coevolution network](#),” *Genome Res*, vol. 19, no. 10, pp. 1861–1871, 2009.
- [33] S. Bradde, A. Braunstein, H. Mahmoudi, F. Tria, M. Weigt, and R. Zecchina, “[Aligning graphs and finding substructures by a cavity approach](#),” *EPL*, vol. 89, no. 3, 2010.
- [34] I. Hajirasouliha, A. Schönhuth, D. de Juan, A. Valencia, and S. C. Sahinalp, “[Mirroring co-evolving trees in the light of their topologies](#),” *Bioinformatics*, vol. 28, no. 9, pp. 1202–1208, 2012.
- [35] M. El-Kebir, T. Marschall, I. Wohlers, M. Patterson, J. Heringa, A. Schönhuth, and G. W. Klau, “[Mapping proteins in the presence of paralogs using units of coevolution](#),” *BMC Bioinformatics*, vol. 14 Suppl 15, p. S18, 2013.
- [36] F. Pazos and A. Valencia, “[Similarity of phylogenetic trees as indicator of protein–protein interaction](#),” *Protein Engineering, Design and Selection*, vol. 14, no. 9, pp. 609–614, 2001.
- [37] D. Juan, F. Pazos, and A. Valencia, “[High-confidence prediction of global interactomes based on genome-wide coevolutionary networks](#),” *Proc. Natl. Acad. Sci. U.S.A.*, vol. 105, no. 3, pp. 934–939, 2008.
- [38] R. Jothi, M. G. Kann, and T. M. Przytycka, “[Predicting protein-protein interaction by searching evolutionary tree automorphism space](#),” *Bioinformatics*, vol. 21 Suppl 1, pp. i241–250, 2005.
- [39] D. Ochoa and F. Pazos, “[Studying the co-evolution of protein families with the Mirrortree web server](#),” *Bioinformatics*, vol. 26, no. 10, pp. 1370–1371, 2010.
- [40] D. Ochoa, D. Juan, A. Valencia, and F. Pazos, “[Detection of significant protein coevolution](#),” *Bioinformatics*, vol. 31, no. 13, pp. 2166–2173, 2015.
- [41] G. Casari, C. Sander, and A. Valencia, “[A method to predict functional residues in proteins](#),” *Nat Struct Mol Biol*, vol. 2, no. 2, pp. 171–178, 1995.
- [42] A. S. Lapedes, B. G. Giraud, L. Liu, and G. D. Stormo, “[Correlated mutations in models of protein sequences: phylogenetic and structural effects](#),” in *Statistics in molecular biology and genetics - IMS Lecture Notes - Monograph Series*, 1999, vol. 33, pp. 236–256.

- [43] F. Morcos, A. Pagnani, B. Lunt, A. Bertolino, D. S. Marks, C. Sander, R. Zecchina, J. N. Onuchic, T. Hwa, and M. Weigt, “[Direct-coupling analysis of residue coevolution captures native contacts across many protein families](#),” *Proc. Natl. Acad. Sci. U.S.A.*, vol. 108, no. 49, pp. E1293–1301, 2011.
- [44] D. S. Marks, L. J. Colwell, R. Sheridan, T. A. Hopf, A. Pagnani, R. Zecchina, and C. Sander, “[Protein 3D structure computed from evolutionary sequence variation](#),” *PLoS ONE*, vol. 6, no. 12, pp. 1–20, 2011.
- [45] R. R. Cheng, F. Morcos, H. Levine, and J. N. Onuchic, “[Toward rationally redesigning bacterial two-component signaling systems using coevolutionary information](#),” *Proc. Natl. Acad. Sci. U.S.A.*, vol. 111, no. 5, pp. E563–571, 2014.
- [46] L. Burger and E. van Nimwegen, “[Accurate prediction of protein-protein interactions from sequence alignments using a Bayesian method](#),” *Mol. Syst. Biol.*, vol. 4, p. 165, 2008.
- [47] A.-F. Bitbol, “[Inferring interaction partners from protein sequences using mutual information](#),” *PLoS Comput. Biol.*, vol. 14, no. 11, p. e1006401, 2018.
- [48] G. Marmier, M. Weigt, and A.-F. Bitbol, “[Phylogenetic correlations can suffice to infer protein partners from sequences](#),” *PLoS Comput. Biol.*, vol. 15, no. 10, p. e1007179, 2019.
- [49] A. Gerardos, N. Dietler, and A.-F. Bitbol, “[Correlations from structure and phylogeny combine constructively in the inference of protein partners from sequences](#),” *PLoS Comput. Biol.*, vol. 18, no. 5, p. e1010147, 2022.
- [50] C. A. Gandarilla-Perez, S. Pinilla, A.-F. Bitbol, and M. Weigt, “[Combining phylogeny and coevolution improves the inference of interaction partners among paralogous proteins](#),” *PLoS Comput. Biol.*, vol. 19, no. 3, p. e1011010, 2023.
- [51] R. M. Rao, J. Liu, R. Verkuil, J. Meier, J. Canny, P. Abbeel, T. Sercu, and A. Rives, “MSA Transformer,” in *Proceedings of the 38th International Conference on Machine Learning*, vol. 139. PMLR, 2021, pp. 8844–8856. [Online]. Available: <https://proceedings.mlr.press/v139/rao21a.html>
- [52] U. Lupo, D. Sgarbossa, and A.-F. Bitbol, “[Protein language models trained on multiple sequence alignments learn phylogenetic relationships](#),” *Nat Commun*, vol. 13, no. 6298, 2022.
- [53] D. Sgarbossa, U. Lupo, and A.-F. Bitbol, “[Generative power of a protein language model trained on multiple sequence alignments](#),” *Elife*, vol. 12, p. e79854, 2023.
- [54] Z. Xie and J. Xu, “[Deep graph learning of inter-protein contacts](#),” *Bioinformatics*, vol. 38, no. 4, pp. 947–953, 2022.
- [55] M. T. Laub and M. Goulian, “[Specificity in two-component signal transduction pathways](#),” *Annu. Rev. Genet.*, vol. 41, pp. 121–145, 2007.
- [56] M. Barakat, P. Ortet, C. Jourlin-Castelli, M. Ansaldi, V. Mejean, and D. E. Whitworth, “[P2CS: a two-component system resource for prokaryotic signal transduction research](#),” *BMC Genomics*, vol. 10, p. 315, 2009.
- [57] M. Barakat, P. Ortet, and D. E. Whitworth, “[P2CS: a database of prokaryotic two-component systems](#),” *Nucleic Acids Research*, vol. 39, no. Database issue, pp. D771–776, 2011.

- [58] A. J. Enright, I. Iliopoulos, N. C. Kyrpides, and C. A. Ouzounis, “Protein interaction maps for complete genomes based on gene fusion events,” *Nature*, vol. 402, no. 6757, pp. 86–90, 1999.
- [59] D. Szklarczyk, R. Kirsch, M. Koutrouli, K. Nastou, F. Mehryary, R. Hachilif, A. L. Gable, T. Fang, N. T. Doncheva, S. Pyysalo, P. Bork, L. J. Jensen, and C. von Mering, “The STRING database in 2023: protein-protein association networks and functional enrichment analyses for any sequenced genome of interest,” *Nucleic Acids Res*, vol. 51, no. D1, pp. D638–D646, 2023.
- [60] T. Paysan-Lafosse, M. Blum, S. Chuguransky, T. Grego, B. L. Pinto, G. A. Salazar, M. L. Bileschi, P. Bork, A. Bridge, L. Colwell, J. Gough, D. H. Haft, I. Letunić, A. Marchler-Bauer, H. Mi, D. A. Natale, C. A. Orengo, A. P. Pandurangan, C. Rivoire, C. J. A. Sigrist, I. Sillitoe, N. Thanki, P. D. Thomas, S. C. E. Tosatto, C. H. Wu, and A. Bateman, “InterPro in 2022,” *Nucleic Acids Research*, vol. 51, no. D1, pp. D418–D427, 2022.
- [61] K. S. Makarova, Y. I. Wolf, S. L. Mekhedov, B. G. Mirkin, and E. V. Koonin, “Ancestral paralogs and pseudoparalogs and their role in the emergence of the eukaryotic cell,” *Nucleic Acids Research*, vol. 33, no. 14, pp. 4626–4638, 2005.
- [62] S. Basu and B. Wallner, “DockQ: A quality measure for protein-protein docking models,” *PLoS ONE*, vol. 11, no. 8, pp. 1–9, 2016.
- [63] P. Bryant and F. Noé, “Improved protein complex prediction with AlphaFold-multimer by denoising the MSA profile,” *bioRxiv*, 2023.
- [64] W. Zheng, Q. Wuyun, and P. L. Freddolino, “Multi-MSA strategy for protein complex structure modeling,” *CASP15 Abstract*, 2022. [Online]. Available: https://predictioncenter.org/casp15/doc/CASP15_Abstracts.pdf
- [65] J. Liu, Z. Guo, T. Wu, R. S. Roy, F. Quadir, C. Chen, and J. Cheng, “Enhancing AlphaFold-Multimer-based protein complex structure prediction with MULTICOM in CASP15,” *bioRxiv*, 2023.
- [66] B. Wallner, “Improved multimer prediction using massive sampling with AlphaFold in CASP15,” *Proteins*, 2023.
- [67] U. Ghani, I. Desta, A. Jindal, O. Khan, G. Jones, N. Hashemi, S. Kotelnikov, D. Padhorny, S. Vajda, and D. Kozakov, “Improved docking of protein models by a combination of Alphafold2 and ClusPro,” *bioRxiv*, 2022.
- [68] K. Olechnovič, L. Valančauskas, J. Dapkunas, and Č. Venclovas, “Prediction of protein assemblies by structure sampling followed by interface-focused scoring,” *bioRxiv*, 2023.
- [69] T. Dao, D. Fu, S. Ermon, A. Rudra, and C. Ré, “FlashAttention: Fast and memory-efficient exact attention with IO-awareness,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 16 344–16 359, 2022. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2022/file/67d57c32e20fd0a7a302cb81d36e40d5-Paper-Conference.pdf
- [70] A. Hawkins-Hooker, D. T. Jones, and B. Paige, “Using domain-domain interactions to probe the limitations of MSA pairing strategies,” in *Machine Learning for Structural Biology Workshop, NeurIPS*, 2022. [Online]. Available: https://www.mlsb.io/papers_2022/Using_domain_domain_interactions_to_probe_the_limitations_of_MSA_pairing_strategies.pdf

- [71] M. Hu, F. Yuan, K. Yang, F. Ju, J. Su, H. Wang, F. Yang, and Q. Ding, “Exploring evolution-aware & -free protein language models as protein function predictors,” in *Advances in Neural Information Processing Systems*, vol. 35, 2022, pp. 38 873–38 884. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2022/hash/fe066022bab2a6c6a3c57032a1623c70-Abstract-Conference.html
- [72] B. Meynard-Piganeau, C. Fabbri, M. Weigt, A. Pagnani, and C. Feinauer, “Generating interacting protein sequences using domain-to-domain translation,” *Bioinformatics*, vol. 39, no. 7, p. btad401, 2023.
- [73] P. Gainza, F. Sverrisson, F. Monti, E. Rodolà, D. Boscaini, M. M. Bronstein, and B. E. Correia, “Deciphering interaction fingerprints from protein molecular surfaces using geometric deep learning,” *Nat Methods*, vol. 17, no. 2, pp. 184–192, 2020.
- [74] J. Tubiana, D. Schneidman-Duhovny, and H. J. Wolfson, “ScanNet: an interpretable geometric deep learning model for structure-based protein binding site prediction,” *Nat Methods*, vol. 19, no. 6, pp. 730–739, 2022.
- [75] L. F. Krapp, L. A. Abriata, F. Cortés Rodríguez, and M. Dal Peraro, “PeSTo: parameter-free geometric deep learning for accurate prediction of protein binding interfaces,” *Nat Commun*, vol. 14, no. 1, p. 2175, 2023.
- [76] M. N. Pun, A. Ivanov, Q. Bellamy, Z. Montague, C. LaMont, P. Bradley, J. Otwinowski, and A. Nourmohammad, “Learning the shape of protein micro-environments with a holographic convolutional neural network,” *bioRxiv*, 2022.
- [77] F. Wu, L. Wu, D. Radev, J. Xu, and S. Z. Li, “Integration of pre-trained protein language models into geometric deep learning networks,” *Commun Biol*, vol. 6, no. 1, p. 876, 2023.
- [78] Y. Si and C. Yan, “Protein language model embedded geometric graphs power inter-protein contact prediction,” *bioRxiv*, 2023.
- [79] J. Su, C. Han, Y. Zhou, J. Shan, X. Zhou, and F. Yuan, “SaProt: Protein language modeling with structure-aware vocabulary,” *bioRxiv*, 2023.
- [80] G. E. Mena, D. Belanger, S. Linderman, and J. Snoek, “Learning latent permutations with Gumbel-Sinkhorn networks,” *6th International Conference on Learning Representations, ICLR 2018 - Conference Track Proceedings*, pp. 1–22, 2018. [Online]. Available: <https://openreview.net/forum?id=Byt3oJ-0W>
- [81] M. D. Zeiler, “ADADELTA: an adaptive learning rate method,” *arXiv*, 2012.
- [82] H. W. Kuhn, “The Hungarian method for the assignment problem,” *Naval Research Logistics Quarterly*, vol. 2, pp. 83–97, 1955.
- [83] L. S. Johnson, S. R. Eddy, and E. Portugaly, “Hidden Markov model speed heuristic and iterative HMM search procedure,” *BMC bioinformatics*, vol. 11, pp. 1–8, 2010.
- [84] M. Steinegger and J. Söding, “MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets,” *Nat Biotechnol*, vol. 35, no. 11, pp. 1026–1028, 2017.
- [85] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirectional transformers for language understanding,” in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, 2019, pp. 4171–4186.

- [86] A. Wang and K. Cho, “BERT has a mouth, and it must speak: BERT as a Markov random field language model,” *arXiv*, 2019.
- [87] K. Goyal, C. Dyer, and T. Berg-Kirkpatrick, “Exposing the implicit energy networks behind masked language models via Metropolis–Hastings,” *arXiv*, 2021.
- [88] R. Rao, J. Meier, T. Sercu, S. Ovchinnikov, and A. Rives, “Transformer protein language models are unsupervised structure learners,” in *International Conference on Learning Representations*, 2021. [Online]. Available: <https://openreview.net/forum?id=fylclEqgvgd>
- [89] C. Norn, B. I. M. Wicky, D. Juergens, S. Liu, D. Kim, D. Tischer, B. Koepnick, I. Anishchenko, F. Players, D. Baker, and S. Ovchinnikov, “Protein sequence design by conformational landscape optimization,” *Proc. Natl. Acad. Sci. U.S.A.*, vol. 118, no. 11, p. e2017228118, 2021.

Supplementary material

S1 Supplementary methods

S1.1 MSA Transformer and masked language modeling for MSAs

We use the MSA Transformer model [51], which takes MSAs as inputs and was trained with a variant of the masked language modeling (MLM) objective [85] on a training set of 26 million MSAs constructed from UniRef50 clusters. The model’s training objective was to correctly predict the identity of randomly masked residue positions in the MSAs in its training set. Specifically, it was trained to minimize an MLM loss, which reads, for an MSA \mathcal{M} , and its masked version $\widetilde{\mathcal{M}}$:

$$\mathcal{L}_{\text{MLM}}(\mathcal{M}, \widetilde{\mathcal{M}}; \theta) = - \sum_{(m,i) \in \text{mask}} \log p(x_{m,i} | \widetilde{\mathcal{M}}; \theta). \quad (\text{S1})$$

Here, $x_{m,i}$ denotes the amino acid at the i -th residue position (column) in the m -th sequence (row) of \mathcal{M} , while θ stands for all the model parameters. At each residue position in the input MSA, MSA Transformer outputs a probability for each of the 21 possible amino-acid and gap symbols, and $p(x_{m,i} | \widetilde{\mathcal{M}}; \theta)$ in Eq. (S1) is the probability associated with the correct residue $x_{m,i}$ at MSA position (m, i) . MSA Transformer’s architecture interleaves multi-headed (tied) row attention blocks and (untied) column attention blocks, over several layers. Therefore, the accessible context for a masked residue consists not only of amino acids at different positions along the same sequence, but also of amino acids from other sequences [51, 52]. This allows the model to capture coevolution information. After pre-training, each term in the right-hand side of Eq. (S1) can be interpreted as the model’s estimate of the (negative) log-likelihood of the amino acid $x_{m,i}$ at a masked position (m, i) [86–88].

S1.2 A differentiable formulation of paralog matching

We formulate a differentiable optimization problem that can be more efficiently solved than the brute-force search, using gradient methods. The goal is to obtain sets of within-species matchings (and thus permutations) that minimize our MLM loss.

The set \mathcal{P}_N of permutation matrices of N objects can be parameterized exactly by square matrices X via the *matching operator*

$$M(X) = \arg \max_{P \in \mathcal{P}_N} [\text{trace}(P^T X)], \quad (\text{S2})$$

which can be computed using standard non-differentiable algorithms for linear assignment problems [82].⁵

We exploit the fact, shown in [80], that permutation matrices can be approximated arbitrarily well by using the *Sinkhorn operator* S , which is defined on square matrices X as follows:

$$S(X) = \lim_{l \rightarrow \infty} S^l(X), \quad \text{where} \quad S^l(X) = (\mathcal{T}_c \circ \mathcal{T}_r)^l(\exp(X)), \quad (\text{S3})$$

\mathcal{T}_c and \mathcal{T}_r are the row- and column-wise normalization operators, and \exp denotes the component-wise matrix exponential.⁶ More precisely, $M(X) = \lim_{\tau \rightarrow 0^+} S(X/\tau)$ for almost all X [80, Theorem 1]. Hence, by choosing a suitably small value of τ , and using S^l [Eq. (S3)] instead of S

⁵More precisely, the right-hand side of Eq. (S2) has a unique solution for almost all X [80].

⁶That is, S^l consists of applying \exp and then iteratively normalizing rows and columns l times.

for a suitably large l , we can define a smooth mapping $\hat{S}(X) = S^l(X/\tau)$ which sends arbitrary square matrices to “soft permutations” approximating *bona fide* (“hard”) permutations. In practice, we use $\tau = 1$ and $l = 10$.

Applying general soft permutations directly on an MSA (after one-hot encoding its residues) yields a dataset consisting of “amino acid mixtures” at each MSA position. Such datasets are out of distribution relative to MSA Transformer’s pre-training since it was trained on single amino acid embeddings, not mixtures of them. Besides, we wish to optimize for an MLM loss defined on realistic MSAs. Therefore, in order to be able to backpropagate through \hat{S} , while also evaluating MLM losses only on MSAs shuffled by hard permutations, we compute the full matching operator M [Eq. (S2)] in the forward pass, but propagate gradients backwards through \hat{S} alone.⁷

S1.3 Datasets

Benchmark prokaryotic datasets. We developed and tested DiffPALM using joint MSAs extracted from two datasets. The first dataset is composed of 23,632 cognate pairs of histidine kinases (HK) and response regulators (RR) from the P2CS database [56, 57], paired using genome proximity, and previously described in [14, 47]. The average number of pairs per species in this dataset is 10.23 (standard deviation: 7.85).

The second dataset consists of 17,950 ABC transporter protein pairs, homologous to the *Escherichia coli* MALG-MALK pair of maltose and maltodextrin transporters, also paired using genome proximity [14, 21]. The average number of pairs per species in this dataset is 5.68 (standard deviation: 5.60). We also considered a similarly constructed dataset of 220 pairs homologous to the *Escherichia coli* NUOA-NUOJ pair of NADH-quinone oxidoreductase subunits, with an average number of pairs per species of 2.04.

Throughout, our focus is on pairing interaction partners among paralogs within each species. In all these benchmark datasets, species comprising only one pair of sequences were discarded. Indeed, pairing is trivial in these cases.

Out of each of these benchmark datasets of known interacting pairs, we consider paired MSAs of depth ~ 50 (resp. ~ 100 or ~ 200), constructed by selecting all the sequences of randomly sampled species from the full dataset. Specifically, for a target MSA depth $\bar{D} = 50, 100$ or 200 , we add randomly sampled complete species one by one; if the first m species (but no fewer) give an MSA depth $D \geq 0.9\bar{D}$, and the first $n \geq m$ species (but no more) give $D \leq 1.1\bar{D}$, then we select the first k species in our final MSA, with k picked uniformly at random between m and n . For such shallow MSAs, existing coevolution-based methods do not perform well. Note also that MSA Transformer’s large memory footprint constrains the depth and length of input MSAs. Concretely, in our GPU (NVIDIA RTX A6000 with 48 GB of memory) it is possible to backpropagate the gradients for an input that has up to $\sim 40,000$ tokens, i.e. ~ 200 sequences of length ~ 200 amino acids.

Eukaryotic complexes. We considered targets whose structures are not in the training set of AFM with v2 weights, and where default AFM predictions are poor. Specifically, we started from those eukaryotic targets from Table A1 of [18] and from the “Benchmark 2” dataset in [7] whose PDB structures were released after the training cutoff for the AFM v2 weights (April 30, 2018). Among those, we focused on multimers with no more than 2 different types of

⁷See [89] for a similar use of “gradient bypassing” in the context of protein design. We write the hard permutation as $[M(X) - \hat{S}(X)] + \hat{S}(X)$, and halt gradient backpropagation through the term in square brackets. Schematically, let L denote the operator which takes an MSA as input, randomly masks it, passes the masked MSA to MSA Transformer, and finally computes the MLM loss. Then, if A is the MSA whose rows we wish to permute, we use $L'(M(X)A)S'(X)A$ instead of $L'(M(X)A)M'(X)A$ as our “gradient”.

monomers, where both monomers come from the same species, and with paired sequences not longer than 500 amino acids, due to GPU memory constraints. Finally, we further restricted to the 15 targets with default AFM predictions yielding the poorest reported DockQ score. They are listed in [Table S1](#). All of them are heterodimers, except 6ABO which is a heterotetramer complex made of two IFFO1 and two XRCC4 molecules. Note that we also show results for AFM with v3 weights (latest release, training cutoff September 30, 2021) in [Fig. S12](#).

S1.4 General points on AFM

For all structure prediction tasks, we use the five pre-trained AFM models with v2 weights [7], except in [Fig. S12](#) where v3 weights are used. We use full genomic databases and code from release v2.3.1 of the official implementation in <https://github.com/deepmind/alphafold>. We use no structural templates, and perform 3 recycles for each structure, without early stopping. We relax all models using AMBER.

When using all pairing methods (default AFM, DiffPALM, orthology-based), we also retained the block MSAs retrieved by the default AFM pipeline [9].

The AFM confidence score is defined as $0.8 \cdot \text{iptm} + 0.2 \cdot \text{ptm}$, where iptm is the predicted TM-score in the interface, and ptm the predicted TM-score of the entire complex [7].

S2 Supplementary tables

PDB ID	L_A	L_B	D	F_{paired}	$\langle d_p \rangle$	MSR	F_{same}	F_{pred}	D_{DiffPALM}	D_{eff}^A	D_{eff}^B
6QU1	322	48	9267	0.02	3.4	5	0.37	0.45	76	20	11
6POG	114	249	18390	0.20	30.1	3	0.03	0.22	821	114	93
6THL	240	185	3515	0.03	2.2	5	0.40	0.46	53	28	26
6L5K	98	113	15857	0.29	22.9	3	0.05	0.75	3478	402	644
6A6I	98	76	4793	0.11	3.1	5	0.30	0.48	259	57	87
5Z5K	380	66	6349	0.03	12	10	0	0.02	3	3	3
6FYH	124	76	15285	0.51	7.1	3	0.14	0.72	5550	1703	1640
6WCW	184	254	20435	0.21	6.7	3	0.18	0.57	2509	263	402
5XLN	190	45	9434	0.04	3.4	5	0.30	0.65	228	30	31
7BQU	114	28	5915	0.04	2.9	5	0.44	0.65	152	67	61
6ABO	227	82	5058	0.10	4.5	3	0.32	0.56	310	65	25
6INE	267	177	19227	0.18	4	3	0.26	0.39	1334	402	480
6IRE	234	194	7599	0.12	4.2	3	0.21	0.63	591	120	210
6PNQ	202	168	5694	0.20	8.7	5	0.13	0.25	282	117	23
6GK2	106	92	3062	0.08	2.4	5	0.35	0.59	145	20	40

Table S1: **Dataset of eukaryotic complexes.** All 15 eukaryotic protein complexes considered here are listed by their PDB ID, and various MSA properties are given. L_A and L_B are the lengths of the aligned amino acid sequences of the two chains A and B considered. D denotes the depth of the full MSA built by AFM, consisting of the paired MSA and the block MSAs (see “[Assessing the impact of pairing on AFM structure prediction](#)”). F_{paired} is the fraction of sequences that are paired by AFM, i.e. the depth of the paired MSA divided by D . $\langle d_p \rangle$ is the average depth of the MSAs of the single species in the AFM-paired MSA. Thus it is the average of the largest depth among the two chains, since the other one is completed by padding sequences of gaps. MSR stands for “maximum size ratio”: if the ratio of the larger to the smaller of the depths of MSAs A and B is larger than MSR, species are not paired by DiffPALM. F_{same} is the fraction of pairs predicted by DiffPALM that is identical to the pairs predicted by the default AFM pairing method. F_{pred} is the ratio of the number of pairs predicted with DiffPALM to the number predicted using the default AFM pairing method. $D_{\text{DiffPALM}} = D \times F_{\text{paired}} \times F_{\text{pred}}$ denotes the number of pairs output by DiffPALM. D_{eff}^A (resp. D_{eff}^B) is the effective depth corrected with phylogenetic weights (with Hamming distance threshold 0.2) [13] of the MSA associated to chain A (resp. chain B) in the MSAs which are paired by DiffPALM. These effective depths quantify MSA diversity. Rows are ordered by increasing mean DockQ score for the default AFM pairing method (cf. [Fig. S5](#)).

MSAs	Pairing method	Pos. Ex.	N_{runs}	Precision-100	Precision-10
HK-RR	Chance	-	-	0.09	-
HK-RR	DCA-IPA [14]	0	-	0.16	-
HK-RR	MI-IPA [47]	0	-	0.15	-
HK-RR	GA-IPA [50]	0	-	0.18	-
HK-RR	ESM-2 (650M)	0	-	0.11 – 0.16	-
HK-RR	DiffPALM-MRA	0	5	0.37	0.67
HK-RR	DiffPALM-MRA	0	20	0.39	0.71
HK-RR	DiffPALM-IPA	0	20 + 10	0.44	-
HK-RR	DiffPALM-MRA	11	5	0.51	0.87
HK-RR	DiffPALM-MRA	19	5	0.61	0.95
HK-RR	DiffPALM-MRA	45	5	0.74	0.99
MALG-MALK	Chance	-	-	0.20	-
MALG-MALK	DCA-IPA [14]	0	-	0.31	-
MALG-MALK	MI-IPA [47]	0	-	0.32	-
MALG-MALK	GA-IPA [50]	0	-	0.42	-
MALG-MALK	ESM-2 (650M)	0	-	0.19 – 0.31	-
MALG-MALK	DiffPALM-MRA	0	5	0.55	0.84

Table S2: Performance of pairing by DiffPALM vs. baselines. We report the pairing precision for variants of DiffPALM, namely MRA & IPA with various numbers of positive examples (Pos. Ex.) and runs (N_{runs}), as well as for different baseline methods, on 40 MSAs comprising about 50 HK-RR or MALG-MALK pairs. With all methods, a full one-to-one within-species pairing is produced for each MSA, and performance is measured by precision, namely, the fraction of correct pairs among predicted pairs, averaged over the 40 MSAs considered. In “precision-100”, this fraction is computed over all predicted pairs (100% of them). In “precision-10”, it is calculated over the top 10% predicted pairs, when ranked by predicted confidence. For the IPA method, we use 20 runs of MRA as starting point, and we add 5 fixed pairs at each new run, see “DiffPALM: Paralog matching based on MLM”. The chance expectation, and the performance of DCA-IPA [14], MI-IPA [47], and GA-IPA [50] are reported as baselines. We also consider a pairing method based on the scores given by the ESM-2 (650M) protein language model [5], see “Pairing based on a single-sequence language model”. For this method, we consider 10 different values of masking probability from 0.1 to 1.0, and we report the range of precisions obtained.

S3 Supplementary figures

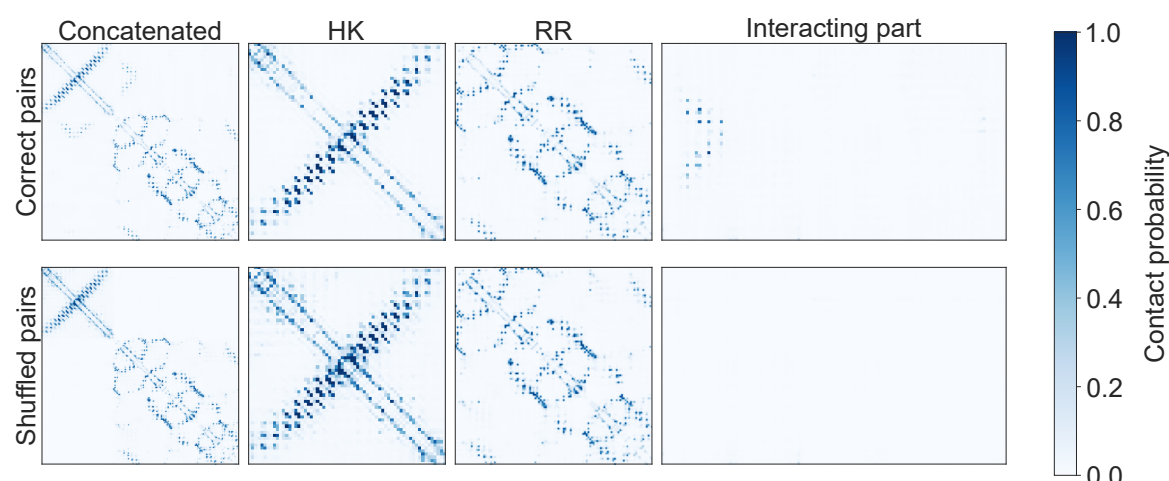


Figure S1: **Comparison of contact maps predicted by MSA Transformer for the correct pairing of an HK MSA and an RR MSA (“Correct pairs”), and for an incorrect pairing (“Shuffled pairs”).** We observe that MSA Transformer is able to correctly predict the inter-protein contacts when given as input a paired MSA made of correctly matched sequences. Conversely, if the model is given as input a paired MSAs where rows have been shuffled before pairing, it is not able to recover the inter-protein contact map (even though it correctly recovers correctly the intra-protein contact maps). These results suggests that MSA Transformer can distinguish between interacting and non-interacting pairs of protein sequences, despite the fact that dimers or paired MSAs were not in the training set used for its MLM pre-training [51].

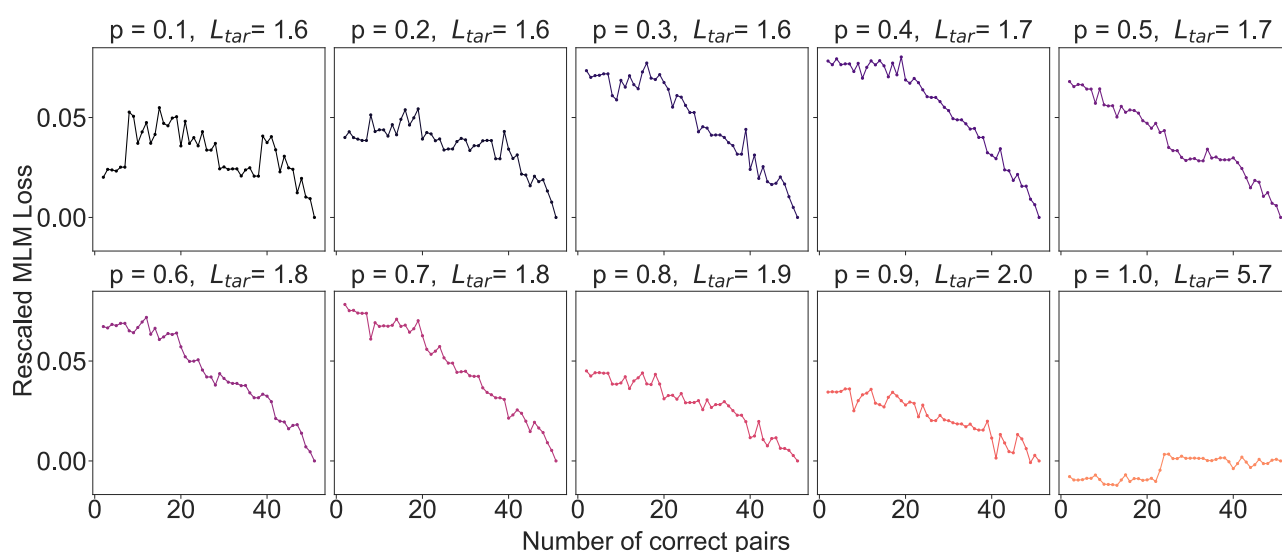


Figure S2: MLM loss vs. number of correct pairs for different masking probabilities. We use an MSA of 50 correctly paired sequences comprising 5 different species from the HK-RR dataset. To estimate the expected loss accurately, we used 20 different masks at each step. L_{tar} denotes the expected loss when all pairs are correctly matched. For visualization purposes, in every plot we rescale the loss by shifting it by L_{tar} . We find that our MLM loss in Eq. (S1) decreases for increasing numbers of correctly matched sequences in the MSA. We see that the sweet spot of the masking probability p (i.e. the value that gives steeper and smoother loss curves) is at moderately high values ($0.4 \leq p \leq 0.8$). As explained in “Methods”, high masking probabilities make it more challenging for the model to predict the masked amino acids using only information coming from the masked MSA, thus encouraging it to use, instead, information coming from the matched MSA. This motivates our choice of a masking probability of $p \geq 0.7$.

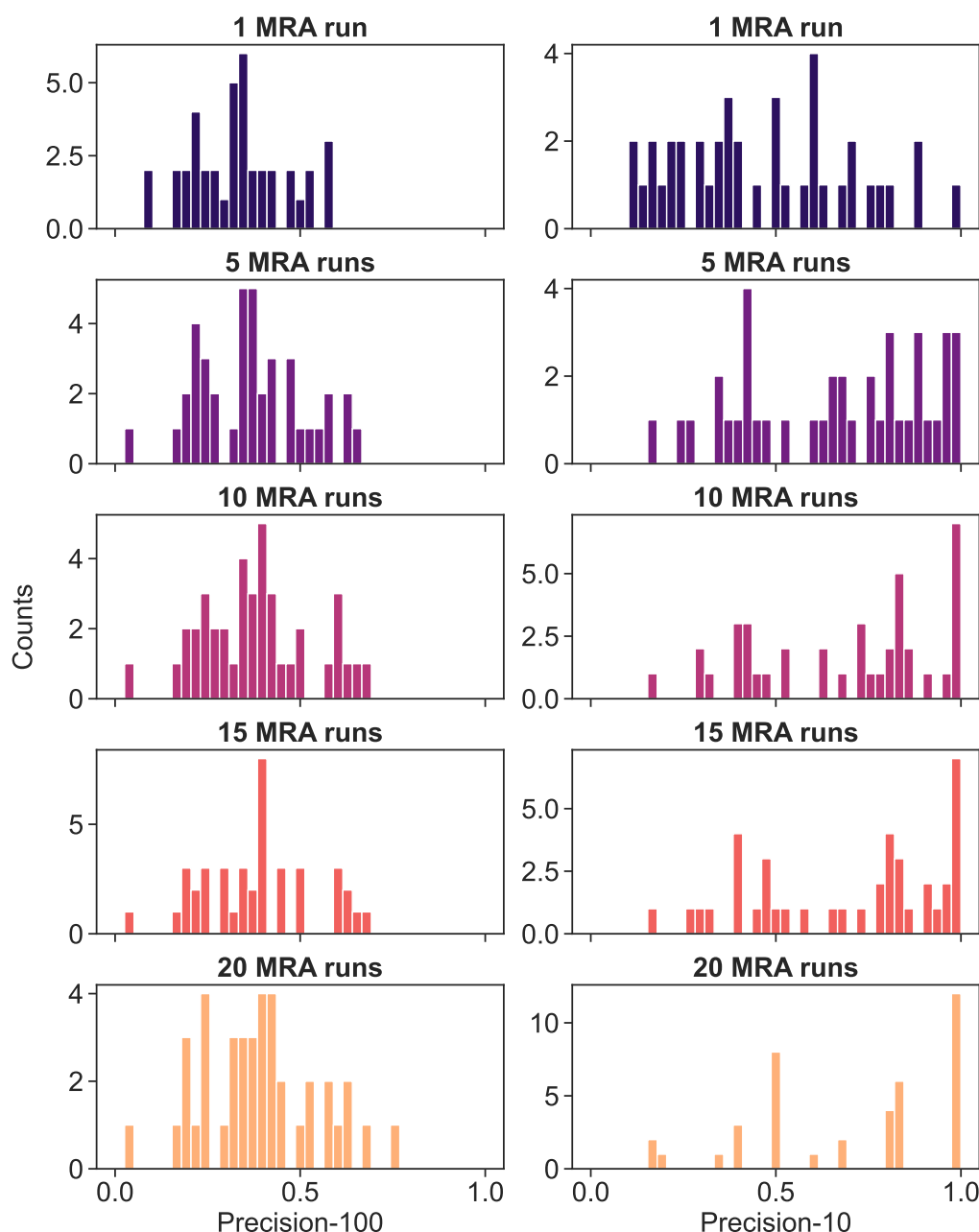


Figure S3: Distribution of precision scores for different number of MRA runs. We report the histograms of precision scores obtained by the MRA variant of DiffPALM on each of the 40 MSAs comprising about 50 HK-RR pairs (used in “[DiffPALM outperforms other coevolution methods on small MSAs](#)”), for different number of MRA runs. Precision-100 and precision-10 are defined in [Fig. 1](#) and [Table S2](#). We observe a skewed distribution for precision-10 scores, especially after many MRA runs: a very high precision is reached for many MSAs, but low precisions are obtained for some. [Fig. 1](#) displays the average and the standard error of each of the distributions shown here.

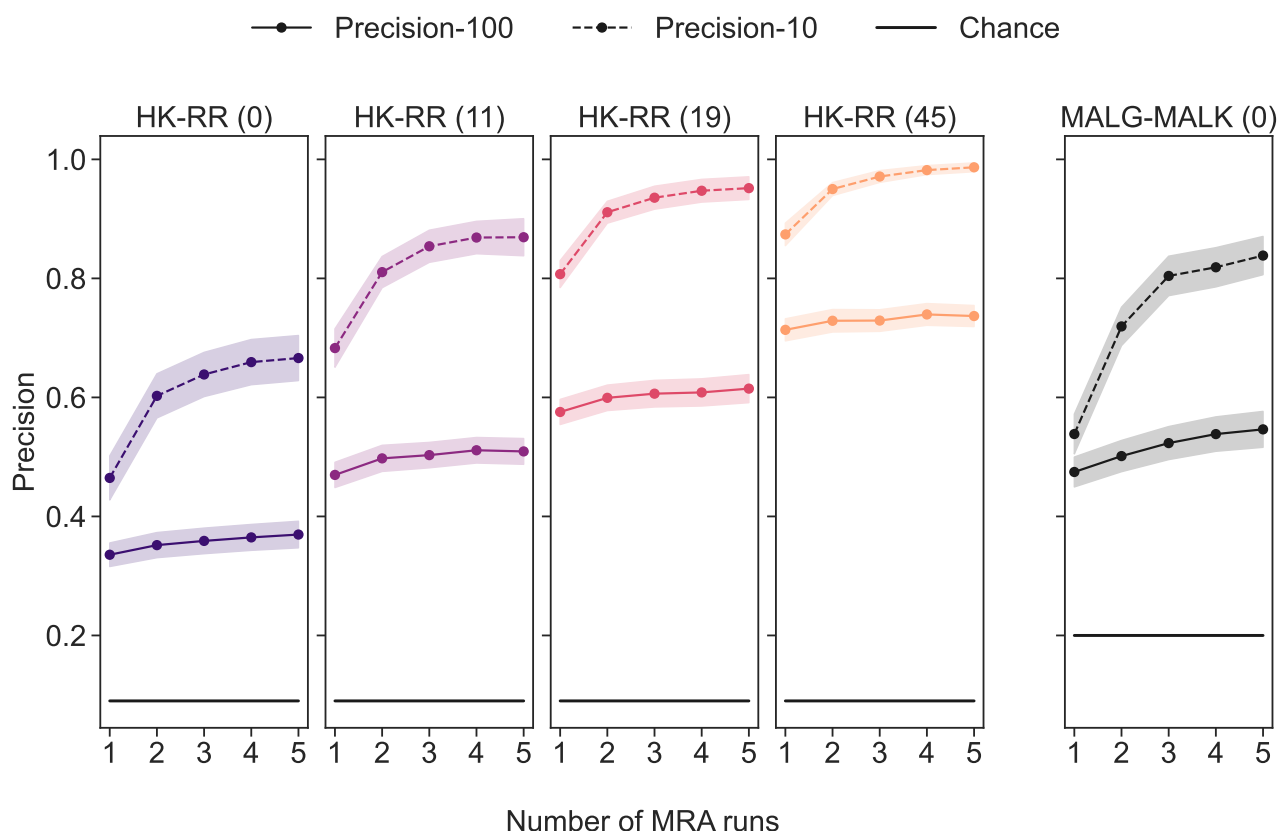


Figure S4: Performance of DiffPALM for different numbers of positive examples and for two pairs of protein families. The performance of DiffPALM is plotted versus the number of MRA runs. As in Fig. 2, we compare runs with various numbers of positive examples (left panels) on the 40 MSAs comprising about 50 HK-RR pairs (used in “DiffPALM outperforms other coevolution methods on small MSAs”), and runs on 40 MSAs comprising about 50 MALG-MALK pairs with no positive example (right panel). Precision-100 and precision-10 are defined in Fig. 1 and Table S2. The protein families considered and the number of positive examples are indicated in the title of each panel (the latter between brackets). In each case, we plot the mean value over the 40 different MSAs considered and the standard error interval. The chance expectation is shown for reference.

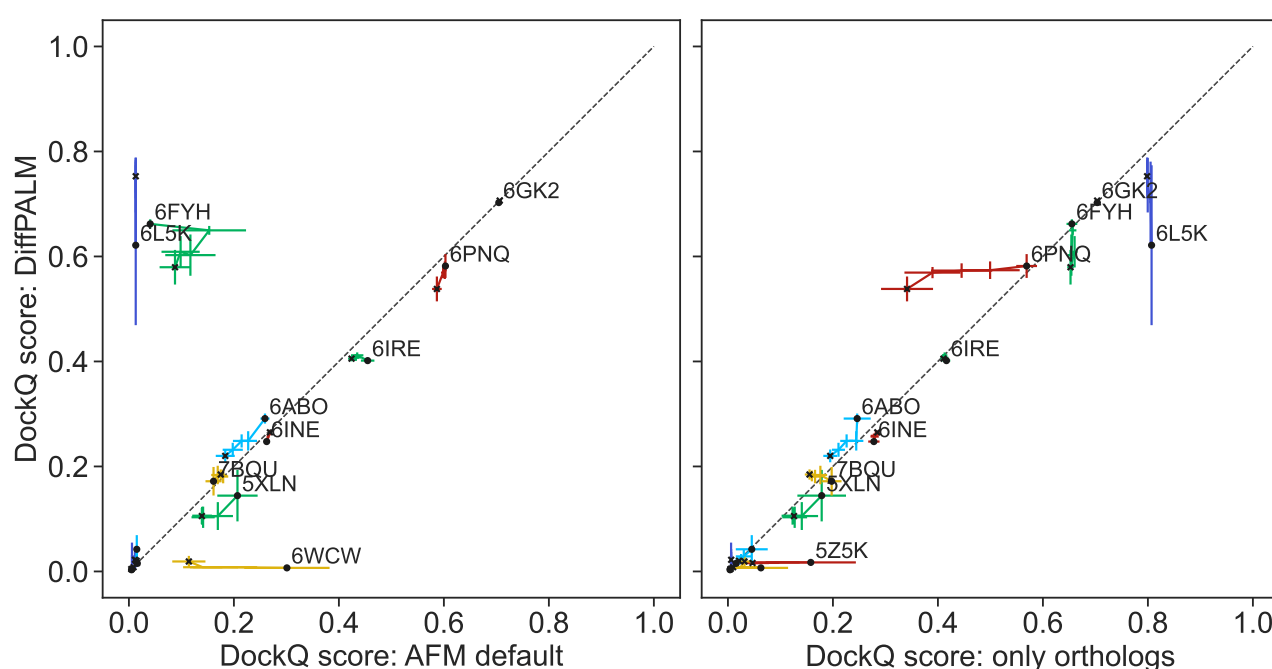


Figure S5: Performance of structure prediction by AFM using different MSA pairing methods. We report the performance of AFM, in terms of DockQ scores, for the 15 complexes listed in Table S1, using three different pairing methods on the same initial unpaired MSAs. Left panel: DiffPALM versus default AFM pairing. Right panel: DiffPALM versus only pairing orthologs to the two query sequences. As in Fig. 3, for each complex, AFM is run five times, and the five top predicted structures by AFM confidence are considered each time, yielding 25 predicted structures total. For each complex, we show “trajectories” of performance starting from the top-confidence predicted structure (black circular marker) and ending with all predicted structures up to and including the fifth one (black cross marker). Results are averaged over the 5 runs and standard errors are shown as error bars. Points with DockQ below 0.1 are not labelled with their PDB ID for graphical reasons. Note that Fig. 3 restricts to those complexes where any two of these three pairing methods yield a significant difference ($> 10\%$) in average DockQ scores, among those shown here and listed in Table S1.

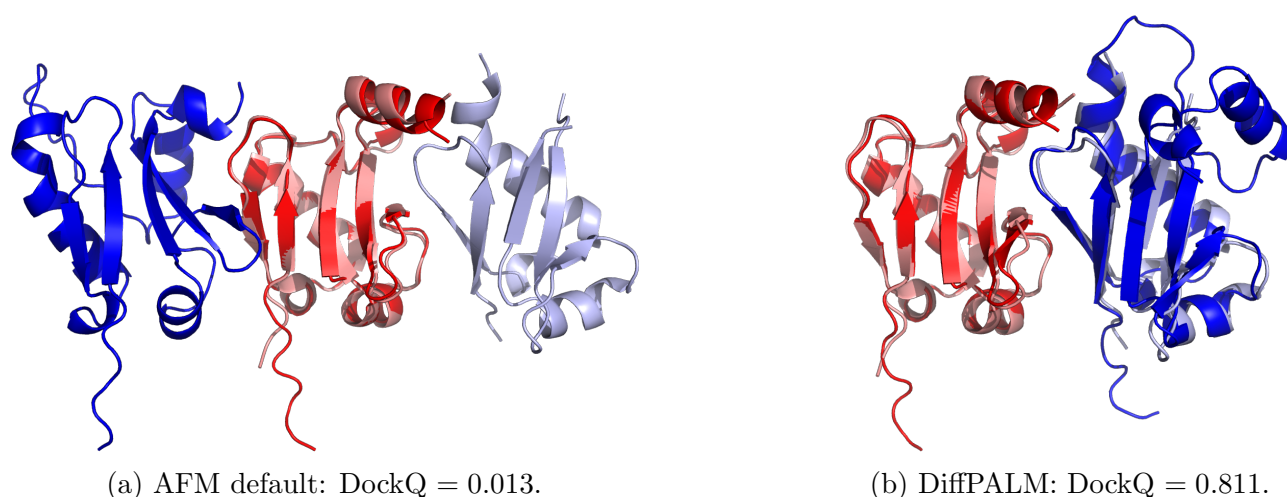


Figure S6: **Comparing the AFM default MSA pairing strategy with DiffPALM, for structure 6L5K.** In both panels, we superimpose the experimental structure of 6L5K with a structure predicted using AFM. Chains A and B of the PDB structure are colored in salmon and light blue respectively, while chains A and B of both predicted structures are colored in bright red and bright blue respectively. (a) Comparing the experimental structure with a typical high-confidence prediction generated with the default MSA pairing pipeline. (b) Comparing the experimental structure with a typical high-confidence prediction generated with our MSA pairing pipeline based on DiffPALM.

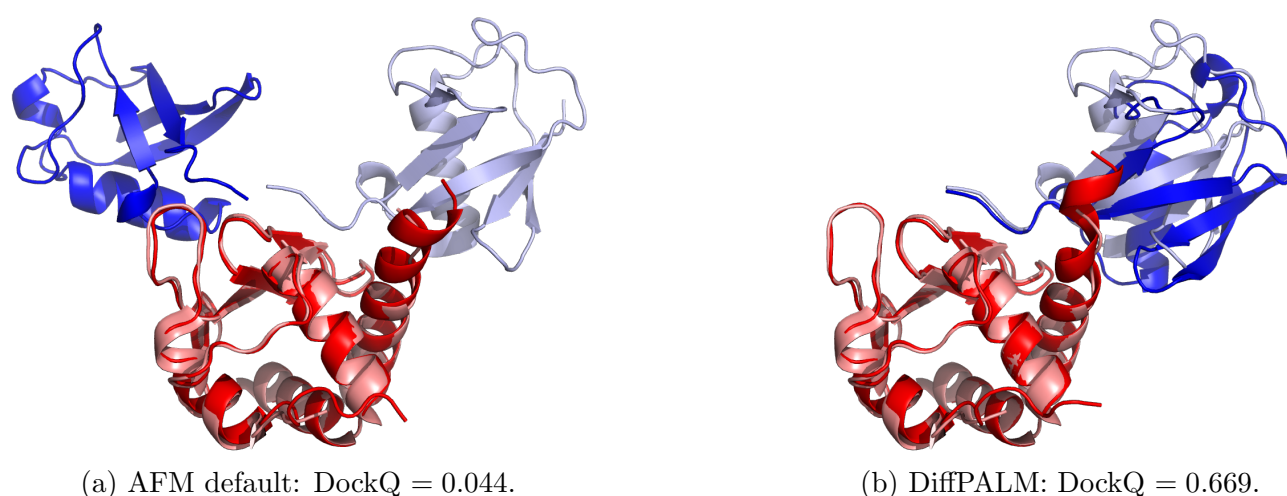


Figure S7: **Comparing the AFM default MSA pairing strategy with DiffPALM, for structure 6FYH.** Same as Fig. S6, but for 6FYH.

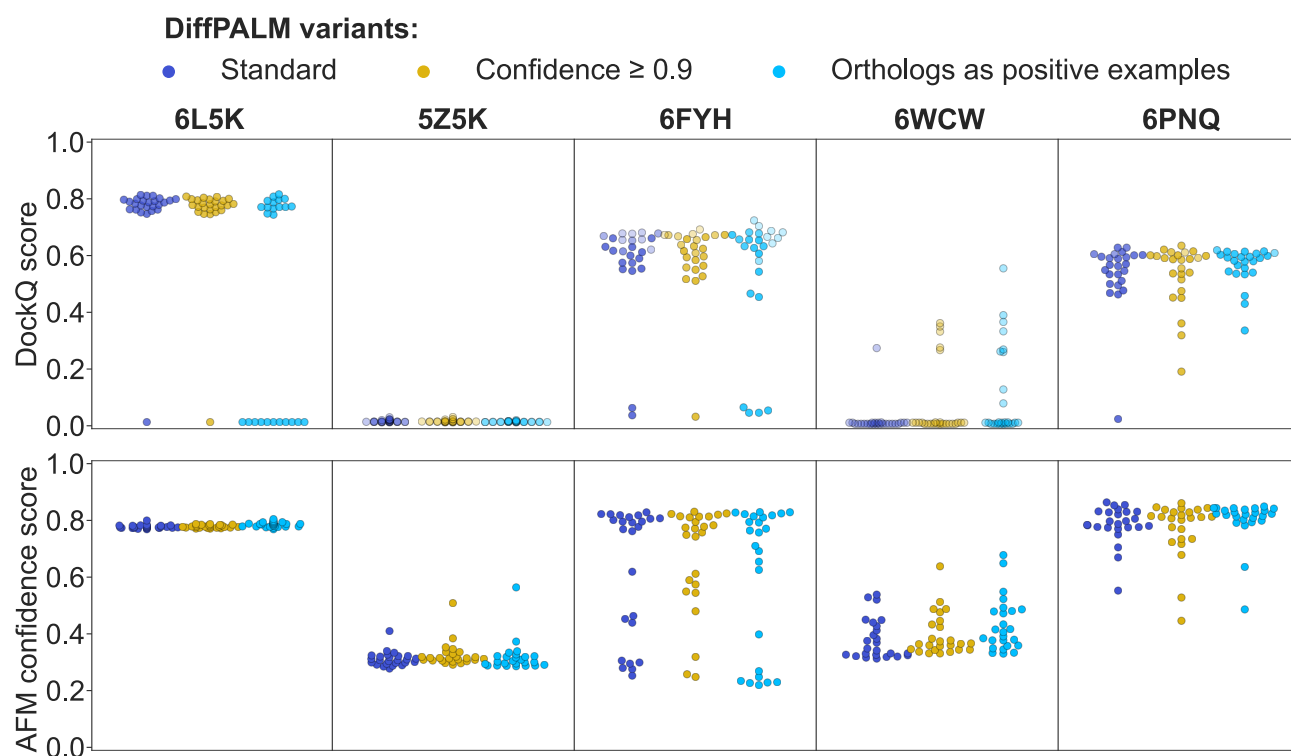


Figure S8: **Performance of AFM using different variants of DiffPALM.** Same as Fig. 3, but here we compare the standard DiffPALM method with two of its variants: one where we only use pairs with high predicted confidence (≥ 0.9) as input to the AFM pipeline, and one where we use orthology-based pairs (i.e. those employed in the “Only orthologs” case shown in Fig. 3) as positive examples for DiffPALM, and use the pairs predicted by DiffPALM, as well as the positive examples, as input of AFM.

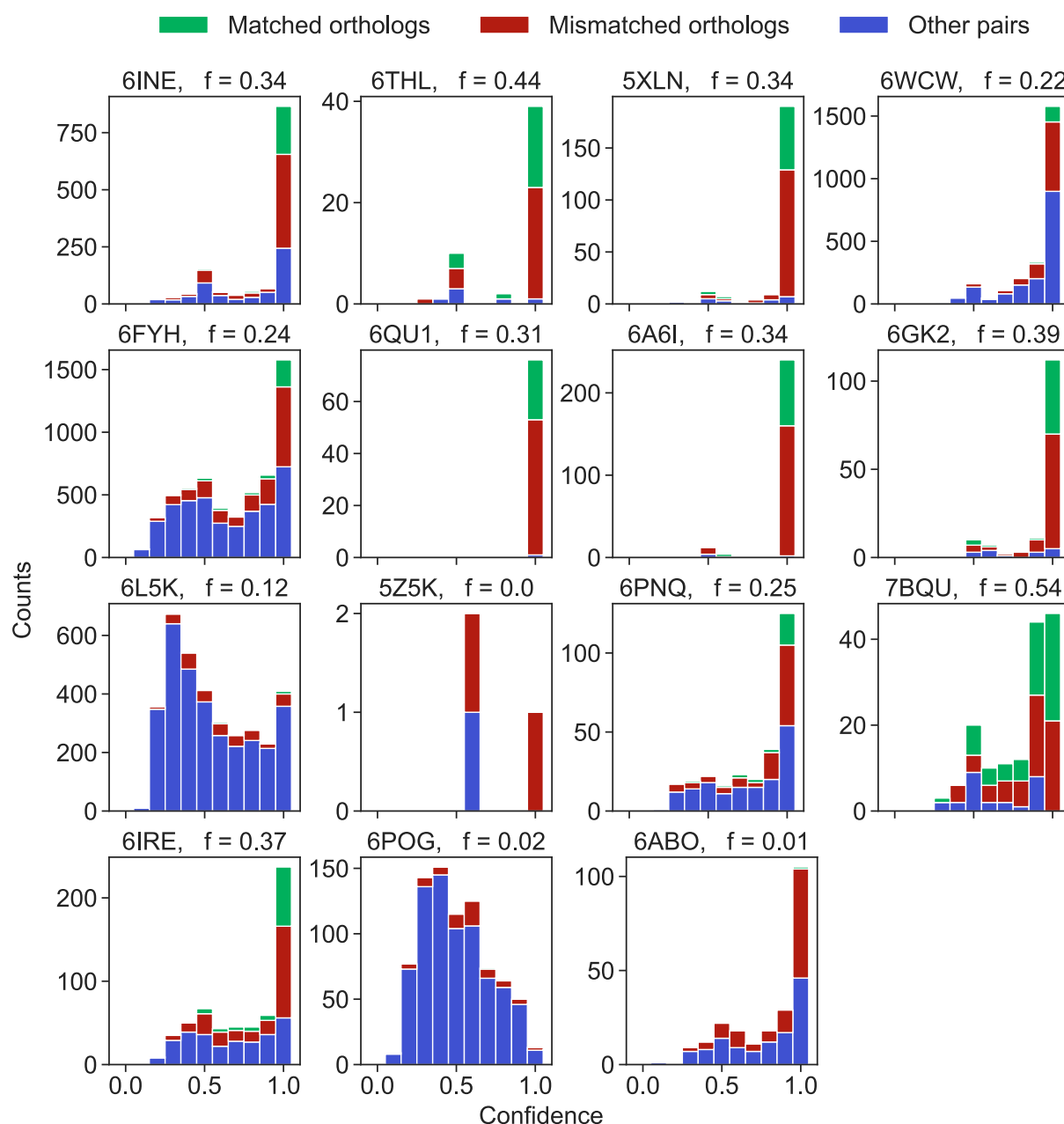


Figure S9: **Confidence of DiffPALM predictions.** We show, for the 15 complexes listed in Table S1, histograms of the DiffPALM confidence values (see “Result and confidence”). We distinguish the orthology-based pairs that are recovered by DiffPALM, the otherwise paired orthologs, and all the other paired sequences. We indicate in panel titles the value of the fraction f of orthology-based pairs that are recovered by DiffPALM.

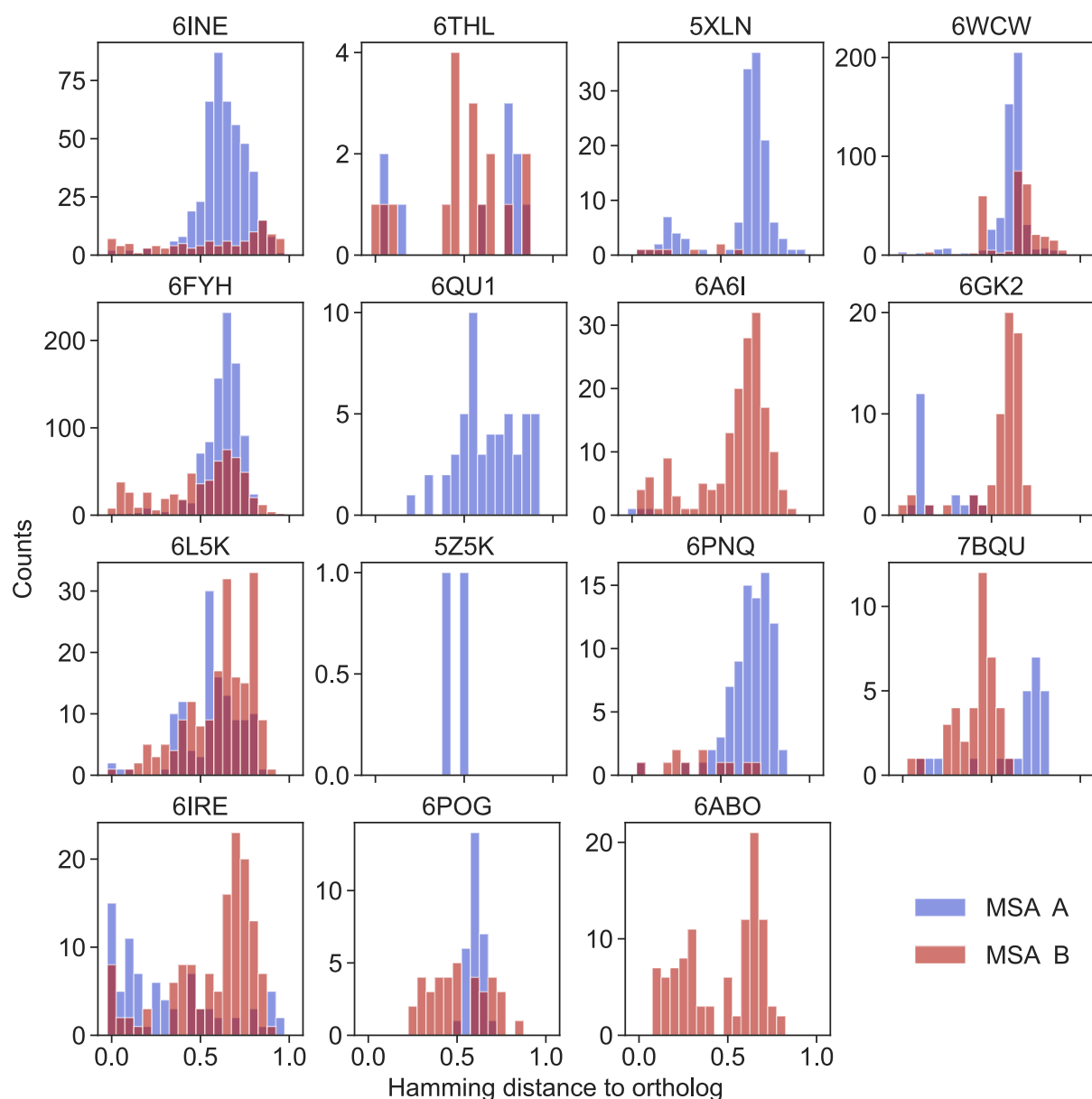


Figure S10: **Hamming distance to the orthologs of all the mismatched pairs predicted by DiffPALM.** We show, for the 15 complexes listed in Table S1, histograms of the Hamming distance between the partner predicted by DiffPALM and the one predicted by matching orthologs to the query sequences, whenever they differ. In practice, for each sequence A_1 in family A which is paired with a partner B_1 from family B using orthology, but with a different partner B_2 using DiffPALM, we measure the Hamming distance between B_1 and B_2 . A similar protocol is conducted for each sequence B_1 in family B. These distances allow us to compare the pairs predicted by DiffPALM to the orthology-based pairs. Note that the total counts of the distributions regarding MSA A and MSA B generally differ. This happens because DiffPALM might pair orthologs to padding sequences of gaps: in this case, we do not report Hamming distances.

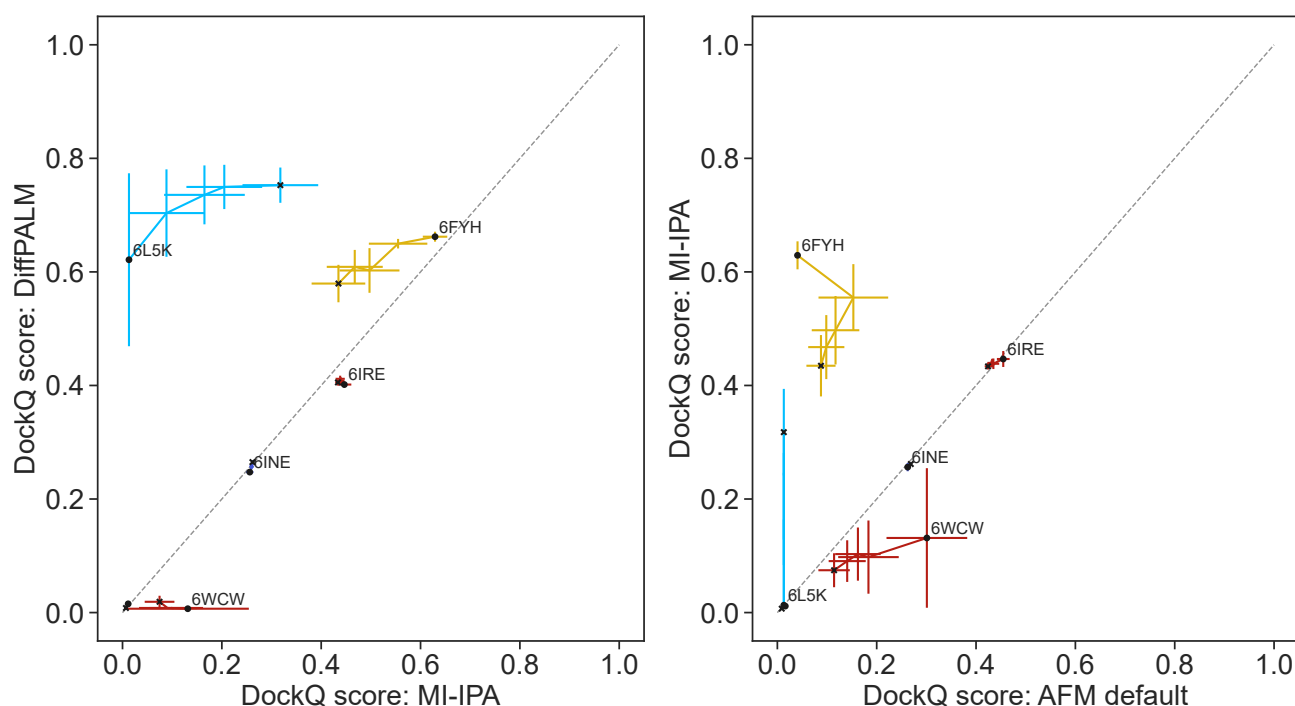


Figure S11: **Performance of structure prediction by AFM using MI-IPA pairing.** Same as Fig. S5 but comparing DiffPALM to MI-IPA (left panel) and MI-IPA to the default AFM pairing method (right panel). Here we restricted to the 6 eukaryotic complexes with deepest pairable MSAs (see Table S1), due to the depth requirements of traditional coevolution methods.

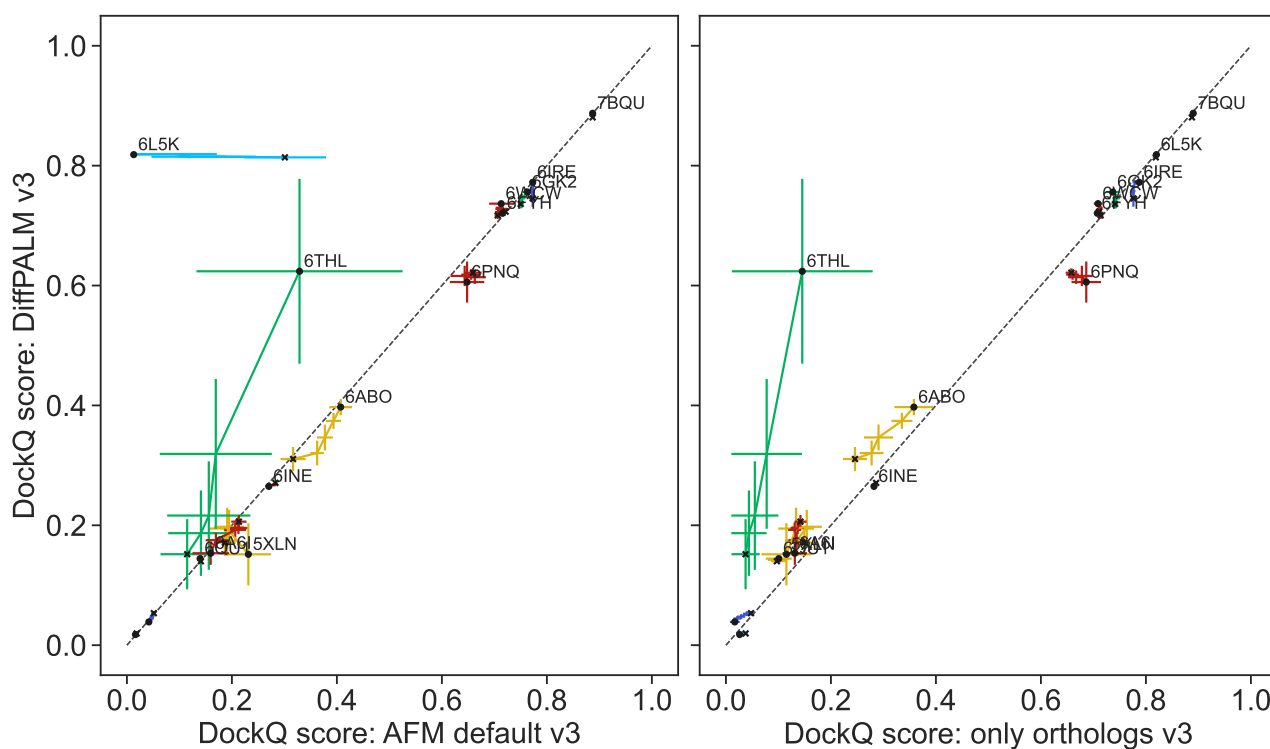


Figure S12: **Performance of structure prediction by AFM-v3 using different MSA pairing methods.** Same as Fig. S5 but using AlphaFold-Multimer with v3 weights (recall that v2 weights are used in the rest of this work).