

ONE PROTEIN IS ALL YOU NEED

Anonymous authors

Paper under double-blind review

ABSTRACT

Generalization beyond training data remains a central challenge in machine learning for biology. A common way to enhance generalization is self-supervised pre-training on large datasets. However, aiming to perform well on all possible proteins can limit a model’s capacity to excel on any specific one, whereas experimentalists typically need accurate predictions for individual proteins they study, often not covered in training data. To address this limitation, we propose a method that enables self-supervised customization of protein language models to one target protein at a time, on the fly, and without assuming any additional data. We show that our Protein Test-Time Training (ProteinTTT) method consistently enhances generalization across different models, their sizes, and datasets. ProteinTTT improves structure prediction for challenging targets, achieves new state-of-the-art results on protein fitness prediction, and enhances function prediction on two tasks. Through two challenging case studies, we also show that customization via ProteinTTT achieves more accurate antibody–antigen loop modeling and enhances 19% of structures in the Big Fantastic Virus Database, delivering improved predictions where general-purpose AlphaFold2 and ESMFold struggle.

1 INTRODUCTION

A comprehensive understanding of protein structure, function, and fitness is essential for advancing research in the life sciences (Subramaniam & Kleywegt, 2022; Tyers & Mann, 2003; Papkou et al., 2023). While machine learning models have shown remarkable potential in protein research, they are typically optimized for achieving the best average performance across large datasets (Jumper et al., 2021; Watson et al., 2023; Kouba et al., 2023). However, biologists often focus their research on individual proteins or protein complexes involved in, for example, metabolic disorders (Ashcroft et al., 2023; Gunn & Neher, 2023), oncogenic signaling (Hoxhaj & Manning, 2020; Keckesova et al., 2017), neurodegeneration (Gulen et al., 2023; oh Seo et al., 2023), and other biological phenomena (Gu et al., 2022). In these scenarios, detailed insights into a single protein can lead to significant scientific advances.

However, general machine learning models for proteins often struggle to generalize to practically interesting individual cases due to data scarcity (Bushuiev et al., 2023; Chen & Gong, 2022) and distribution shifts (Škrinjar et al., 2025; Tagasovska et al., 2024; Feng et al., 2024). Bridging the gap between broad, dataset-wide optimization and precision needed to study single proteins of practical interest remains a key challenge in integrating machine learning into biological research (Sapoval et al., 2022). This challenge is particularly acute in computational biology, where accurate predictions for individual proteins are essential to guide resource-intensive wet-lab experiments, in contrast to

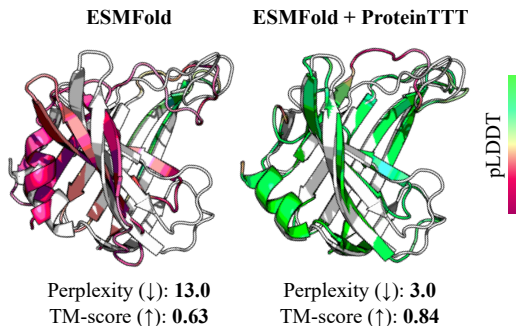


Figure 1: **Example of protein structure prediction after single-protein model customization via ProteinTTT.** ESMFold poorly predicts the structure of the CASP14 target T1074 (white) because the underlying language model ESM2 poorly fits the sequence, as indicated by the high perplexity (left and Fig. 2E in Lin et al. (2023)). Self-supervised test-time customization of ESM2 to the single sequence of T1074 reduces the perplexity, resulting in improved structure prediction (right).

domains such as natural language processing or computer vision, where models are typically expected to flexibly handle diverse prompts from many users in real time (Brown, 2020; Ramesh et al., 2021).

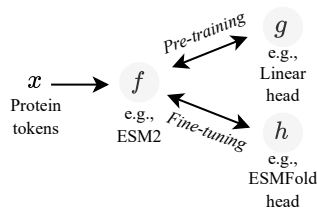
To address this challenge, we propose a test-time approach for generalization to one protein at a time, effectively enabling more accurate predictions for individual targets, particularly those poorly represented in training data. Our Protein Test-Time Training (ProteinTTT) method customizes protein language models (PLMs) to individual proteins on the fly and without assuming additional data. Our approach is based on a simple yet powerful premise: if a language model is less perplexed (surprised) by a protein sequence—or if it “understands” its unique patterns better—it will generate a more accurate representation for predicting its structure and function. Given a model pre-trained via masked language modeling, our method effectively minimizes perplexity on a target protein or its multiple sequence alignment (MSA) through self-supervised customization, improving downstream performance without updating the downstream task head. The widespread use of masked modeling as a pre-training paradigm makes ProteinTTT broadly applicable in computational biology.

In summary, this work demonstrates the surprising effectiveness of protein model customization and lays the foundation for exploring other test-time strategies and broader biological applications. The key contributions are: **(1)** We introduce ProteinTTT, to the best of our knowledge the first customization method in machine learning for biology. We provide a user-friendly and easily extensible implementation¹ and provide insights into the effectiveness of protein model customization by linking it to perplexity minimization. **(2)** We empirically validate ProteinTTT, showing improvements in protein structure prediction with well-established models, achieving state-of-the-art results in protein fitness prediction, and enhancing protein function prediction on terpene synthase substrate classification and protein localization prediction. **(3)** We demonstrate the practical utility of focusing on one protein at a time through two challenging case studies. ProteinTTT enables more accurate prediction of antibody–antigen loops and improves 19% of structures in the Big Fantastic Virus Database, delivering accurate predictions where general-purpose AlphaFold2 and ESMFold struggle.

2 BACKGROUND AND RELATED WORK

The broad adoption of Y-shaped architectures relying on masked modeling enables the development of a general method for customizing protein models at test time via masking-based self-supervision.

The Y-shaped paradigm of learning. In machine learning applied to proteins, architectures often follow a Y-shaped paradigm (Gandelsman et al., 2022), consisting of a backbone feature extractor f operating on protein tokens x , a self-supervised head g , and an alternative fine-tuning head h . During training, $g \circ f$ is first pre-trained, and the pre-trained backbone f is then reused to fine-tune $h \circ f$ toward a downstream task. Here, \circ denotes a composition of two machine learning modules (e.g., g is applied on top of f in $g \circ f$). At test time, the final model $h \circ f$ is fixed. Generalization is achieved by leveraging the rich knowledge encoded in the backbone f and the task-specific priors embedded in the fine-tuning head h . This paradigm enables overcoming data scarcity during fine-tuning and underlies breakthrough approaches in protein structure prediction (Lin et al., 2023), protein design (Watson et al., 2023), protein function prediction (Yu et al., 2023), and other tasks (Hayes et al., 2024).



The backbone f is typically a large neural network pre-trained in a self-supervised way on a large dataset using a smaller pre-training projection head g (Hayes et al., 2024). The fine-tuning head h , however, depends on the application. In some cases, h is a large neural network, repurposing the pre-trained model entirely (Watson et al., 2023; Lin et al., 2023); in others, h is a minimal projection with few parameters (Cheng et al., 2023), or even without any parameters at all (i.e., a zero-shot setup; Meier et al. (2021); Dutton et al. (2024)). The fine-tuning head h can also be a machine learning algorithm other than a neural network (Samusevich et al., 2025).

Masked modeling. While the objective of fine-tuning $h \circ f$ is determined by the downstream application, the choice of pre-training objective for $g \circ f$ is less straightforward. Nevertheless, the

¹<https://anonymous.4open.science/r/ProteinTTT-anonymous-F585>

108 dominant paradigm for protein pre-training is masked modeling, which optimizes model weights to
 109 reconstruct missing protein parts. This objective has proven effective across diverse tasks (Heinzinger
 110 & Rost, 2025; Schmirler et al., 2024), including structure (Lin et al., 2023; Jumper et al., 2021),
 111 fitness (Meier et al., 2021; Su et al., 2023), and function prediction (Samusevich et al., 2025; Yu
 112 et al., 2023; Elnaggar et al., 2021), as well as protein design (Hsieh et al., 2025; Hayes et al., 2024;
 113 Nijkamp et al., 2023), and has been successfully applied to various protein representations such as
 114 sequences (Hayes et al., 2024; Elnaggar et al., 2023), graphs (Dieckhaus et al., 2024; Bushuiev et al.,
 115 2023), and voxels (Diaz et al., 2023).

116 **Model customization.** Several studies have shown that machine learning models for proteins benefit
 117 from being fine-tuned on protein-specific (Notin et al., 2024; Kirjner et al., 2023; Rao et al., 2019)
 118 or protein family-specific (Sevgen et al., 2023; Samusevich et al., 2025) data. However, collecting
 119 additional data may be resource-intensive, and for many targets, relevant datasets or proteins may be
 120 limited or not available (Durairaj et al., 2023; Kim et al., 2025). In this paper, we propose a versatile
 121 method enabling customizing PLMs for a single target protein or its MSA in a self-supervised manner,
 122 on the fly, and without assuming any additional data. Customization methods have been developed in
 123 computer vision (Chi et al., 2024; Wang et al., 2023; Xiao et al., 2022; Karani et al., 2021) and natural
 124 language processing (Hübötter et al., 2024; Hardt & Sun, 2023; Ben-David et al., 2022; Banerjee
 125 et al., 2021). The paradigm of test-time training (TTT), developed to mitigate distribution shifts in
 126 computer vision applications (Gandelsman et al., 2022; Sun et al., 2020), is the main inspiration for
 127 our work. We demonstrate that customization via test-time training enhances the accuracy of PLMs
 128 across a wide range of downstream tasks even without the presence of explicit distribution shifts.

130 3 PROTEIN MODEL CUSTOMIZATION WITH PROTEINTTT

131
 132 In this section, we describe the proposed Protein Test-Time Training (ProteinTTT) approach (Sec-
 133 tion 3.1), followed by its applications to a range of well-established models and datasets (Section 3.2).

135 3.1 SELF-SUPERVISED CUSTOMIZATION TO A TARGET PROTEIN

136
 137 At test time, we assume a Y-shaped model with a backbone f that has been pre-trained via the
 138 self-supervised track $g \circ f$, followed by task-specific fine-tuning through the supervised track $h \circ f$.
 139 The goal of customization with ProteinTTT is to adapt the backbone f to a single protein x before
 140 making a prediction on a downstream task via the supervised track $h \circ f$. To achieve this, we
 141 customize the backbone f to the single example x :

$$142 \text{ProteinTTT} : (h \circ f(\cdot; \theta_0), x) \mapsto h \circ f(\cdot; \theta_x) \quad (1)$$

143 where θ_0 denotes pre-trained parameters and θ_x parameters optimized for the target protein x using
 144 the self-supervised track $g \circ f$, while the supervised head h remains frozen. Figure 2a illustrates our
 145 customization approach, which is summarized in the following sections. Appendix C describes the
 146 extension of our method to customization using a MSA of a protein, rather than its single sequence.

147
 148 **Customization training objective.** We customize $g \circ f$ to a single target protein sequence x via
 149 minimizing the masked language modeling objective (Devlin, 2018; Rives et al., 2021):

$$151 \mathcal{L}(x; \theta) = \mathbb{E}_{M \sim p_{\text{mask}}(M)} \left[\sum_{i \in M} -\log p(x_i | x_{\setminus M}; \theta) \right], \quad (2)$$

152
 153 where x denotes a sequence of protein tokens (typically amino acid types), and \mathbb{E}_M represents the
 154 expectation over randomly sampled masking positions M . The objective function $\mathcal{L}(x; \theta)$ maximizes
 155 the log-probabilities $\log p(x_i | x_{\setminus M}; \theta) \doteq g(f(x_{\setminus M}; \theta))_i$ of the true (i.e., wild-type) tokens x_i at the
 156 masked positions $i \in M$ in the partially masked sequence $x_{\setminus M}$, where θ denotes the parameters of the
 157 backbone f , and g is the masked language modeling head. **While we focus on classical bi-directional
 158 masked modeling, we also demonstrate that ProteinTTT can be similarly applied to autoregressive
 159 and discrete diffusion models (Appendix B).**

160
 161 To ensure consistency between the customization and pre-training, ProteinTTT adopts the same
 masking and data preprocessing strategies used during pre-training. Specifically, $p_{\text{mask}}(M)$ can

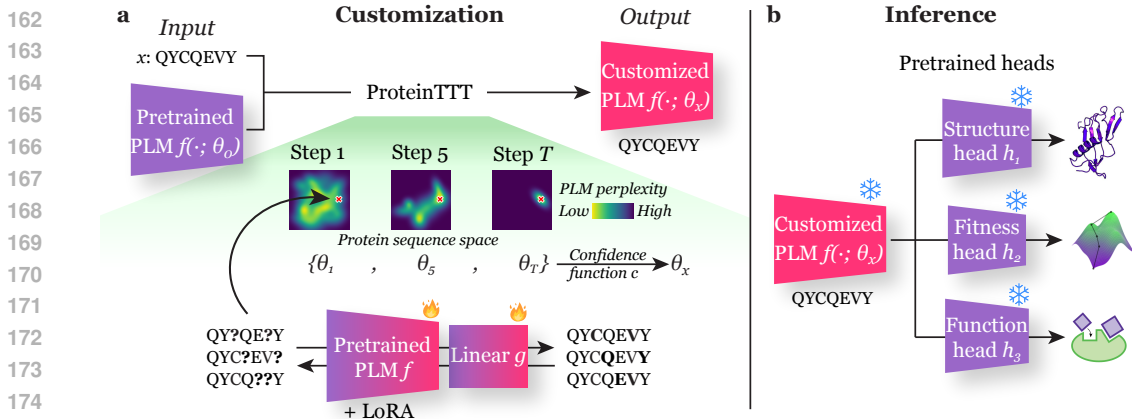


Figure 2: **Overview of protein language model (PLM) customization with ProteinTTT.** (a) Given a protein sequence of interest x and a pretrained PLM $f(\cdot; \theta_0)$, ProteinTTT yields a customized version of the PLM $f(\cdot; \theta_x)$ for that sequence. Customization is achieved by fine-tuning (fire icon) the pretrained parameters θ_0 via masked language modeling solely on the input sequence for T steps, selecting the optimal parameters θ_x using a confidence function c . This procedure adapts the model specifically to the input sequence, improving its internal representation as measured by model perplexity. (b) Once customized, the PLM can be used with pretrained task-specific heads, such as structure, fitness, or function prediction modules, h_1 , h_2 , and h_3 , respectively, without modifying their parameters (snowflake icon). For example, the ESM2 PLM can be customized and then used with the pretrained ESMFold structure prediction head without modifying its 1.4-billion task-specific parameters, resulting in improved structure prediction for the given sequence (e.g., Figure 1).

follow different distributions, such as sampling a fixed proportion (e.g., 15%) of random amino acid tokens (Lin et al., 2023), or dynamically varying the number of sampled tokens based on another distribution (e.g., a beta distribution; Hayes et al. (2024)). During the customization, we replicate the masking distribution used during the pre-training. We also replicate other pre-training practices, such as replacing 10% of masked tokens with random tokens and another 10% with the original tokens (Devlin, 2018; Lin et al., 2023; Su et al., 2023) or cropping sequences to random 1024-token fragments (Lin et al., 2023; Su et al., 2023).

Optimization. Since customization with ProteinTTT does not assume more than a single protein, early stopping on validation data is not feasible. To address this, we first fine-tune the pre-trained parameters θ_0 of a backbone f for a fixed number of steps T , yielding parameters $\Theta = \{\theta_0, \theta_1, \dots, \theta_T\}$. The final customized parameters θ_x are selected as $\arg \max_{\theta \in \Theta} c(h(f(x; \theta)))$ where c is a confidence function. If c is not available, we set $\theta_x = \theta_T$. Appendix H.2 discusses how using pLDDT as the confidence function c for structure prediction makes ProteinTTT robust to hyperparameter selection and how the number of steps T can be fixed (e.g., $T = 30$) while optimizing learning rate and batch size effectively. Before customizing for the next target protein, the parameters are reset to θ_0 .

To make ProteinTTT easily applicable to large-scale models (e.g., the 3B-parameter ESM2 backbone), we leverage low-rank adaptation (LoRA; Hu et al. (2021)) and gradient accumulation during customization. Additionally, to improve the stability and predictability of customization, we use stochastic gradient descent (SGD; Ruder (2016)) instead of the commonly used Adam optimizer (Kingma & Ba, 2015), following (Gandelsman et al., 2022). Further details are provided in Appendix F.

3.2 INFERENCE ON DOWNSTREAM TASKS

Once the backbone f is adapted to a target protein via self-supervised customization, it can be used in conjunction with a pre-trained downstream head h , as $h \circ f$. The key idea of customization with ProteinTTT is not to update the head h , but instead to leverage improved representations from f (Figure 2b). Appendix A provides a justification for why these customized representations generally enhance performance on downstream tasks by linking ProteinTTT to perplexity minimization.

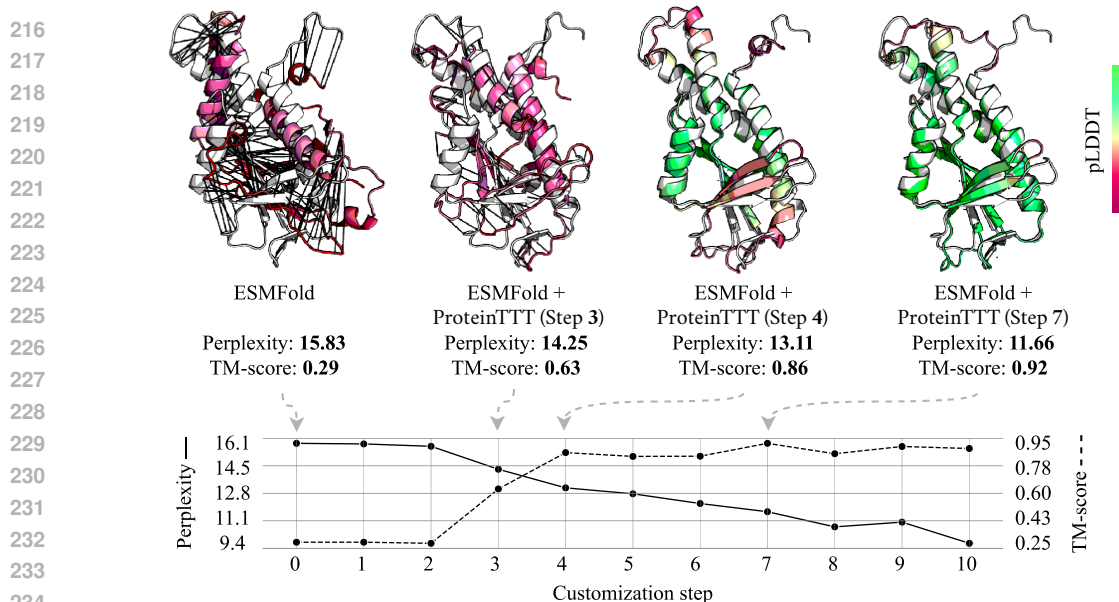


Figure 3: **Customization with ProteinTTT improves protein structure prediction by reducing protein sequence perplexity.** ESMFold fails to predict the structure of chain B from PDB entry 7EBL in the CAMEO validation set, as shown at customization step 0, where the perplexity is high and the TM-score is low. By applying customization with ProteinTTT for the single target sequence, the model iteratively improves the structure prediction quality, as demonstrated by the increasing TM-score, associated with reduced perplexity. At customization step 7, the predicted structure achieves the highest TM-score, as well as the highest predicted confidence metric pLDDT, enabling the selection of this step as the final prediction by the customized ESMFold + ProteinTTT.

Since Y-shaped architectures are prevalent in protein machine learning, ProteinTTT can be straightforwardly applied to numerous tasks. In this work, we consider three standard problems: protein structure, fitness, and function prediction, and apply our method to corresponding well-established models. For structure prediction, we apply ProteinTTT to ESMFold (Figure 3, Lin et al. (2023)), HelixFold-Single (Fang et al., 2023), DPLM2 Bit-based (Hsieh et al., 2025), and ESM3 (Hayes et al., 2024); for fitness prediction, we use ESM2 (Lin et al., 2023), SaProt (Su et al., 2023), ProSST (Li et al., 2024), MSA Transformer (Rao et al., 2021), and ProGen2 (Nijkamp et al., 2023); and for function prediction, we apply ProteinTTT to ESM-1v-based (Meier et al., 2021) EnzymeExplorer (Samusevich et al., 2025) and ESM-1b-based (Rives et al., 2021) Light attention (Stärk et al., 2021).

In all models we consider, f is a Transformer encoder operating on protein tokens, and g is a masked language modeling head mapping embeddings to amino acid types. The downstream head h , however, varies strongly by task. For structure prediction, h is a structure predictor: AlphaFold2-inspired modules in ESMFold, HelixFold-Single and DPLM2 Bit-wise (Jumper et al., 2021), and a VQ-VAE decoder in ESM3 (Razavi et al., 2019). For fitness prediction, h outputs a single score; all methods perform zero-shot inference using $h \circ f$ via log likelihoods from g , with h acting as a simple, parameter-free adaptation of g . For function prediction, h is a classifier: a random forest in EnzymeExplorer (Samusevich et al., 2025) and a light attention module in (Stärk et al., 2021).

4 EXPERIMENTS

In this section, we evaluate ProteinTTT on three well-established downstream tasks in protein machine learning: structure (Section 4.1), fitness (Section 4.2), and function (Section 4.3) prediction.

4.1 PROTEIN STRUCTURE PREDICTION

Protein structure prediction is the task of predicting 3D atom coordinates from an amino acid sequence. It is arguably one of the best-established problems in computational biology (Jumper et al., 2021).

Table 1: **Customization with ProteinTTT improves protein structure prediction.** The metrics are averaged across 18 ESMFold low-confidence targets in the CAMEO test set, and standard deviations correspond to 5 random seeds. CoT and MP stand for the chain of thought and masked prediction baselines.

Method	TM-score \uparrow	LDDT \uparrow
ESM3 (Hayes et al., 2024)	0.3480 \pm 0.0057	0.3723 \pm 0.0055
ESM3 + CoT (Hayes et al., 2024)	0.3677 \pm 0.0088	0.3835 \pm 0.0024
ESM3 + ProteinTTT (Ours)	0.3954 \pm 0.0067	0.4214 \pm 0.0054
DPLM2 Bit-based (Hsieh et al., 2025)	0.3701 \pm 0.0102	0.4681 \pm 0.0071
DPLM2 Bit-based + ProteinTTT (Ours)	0.3796 \pm 0.0024	0.4742 \pm 0.0093
HelixFold-Single (Fang et al., 2023)	0.4709	0.4758
HelixFold-Single + ProteinTTT (Ours)	0.4839 \pm 0.0045	0.4840 \pm 0.0061
ESMFold (Lin et al., 2023)	0.4649	0.5194
ESMFold + MP (Lin et al., 2023)	0.4862 \pm 0.0043	0.5375 \pm 0.0070
ESMFold + ProteinTTT (Ours)	0.5047 \pm 0.0132	0.5478 \pm 0.0058

Evaluation setup. To evaluate the performance of ProteinTTT, we employ CAMEO, a standard benchmark for protein folding. We use the validation and test folds from Lin et al. (2023), focusing only on targets with low-confidence predictions from the base ESMFold, as determined by pLDDT and perplexity (Appendix F.1). We use the standard TM-score (Zhang & Skolnick, 2004) and LDDT (Mariani et al., 2013) metrics to evaluate global and local structure prediction quality, respectively.

As baseline methods, we use techniques alternative to ProteinTTT for improving the performance of the pre-trained base models. In particular, the ESMFold paper proposes randomly masking 15% of amino acids in a protein sequence before inference, allowing for sampling multiple protein structure predictions from the regression ESMFold model (Lin et al., 2023). For each sequence, we sample a number of predictions equal to the total number of ProteinTTT steps and refer to this baseline as ESMFold + MP (Masked Prediction). As a baseline for ESM3, we use chain-of-thought iterative decoding, referred to as ESM3 + CoT, proposed in the ESM3 paper (Hayes et al., 2024).

Results. Customization with ProteinTTT consistently improves the performance of all the tested methods, ESMFold, HelixFold-Single, and ESM3, outperforming the masked prediction (ESMFold + MP) and chain-of-thought (ESM3 + CoT) baselines, as shown in Table 1. Among the 18 challenging CAMEO test proteins, ProteinTTT significantly improved the prediction of 7, 4, 5, and 6 structures from ESMFold, DPLM2 Bit-based, HelixFold-Single, and ESM3, respectively, while only moderately disrupting the prediction of 2, 1, 1, and 1 structures, respectively (Figure A6). Remarkably, ProteinTTT improves DPLM2 Bit-based despite the absence of a confidence function (no trained pLDDT head available) and despite the model being pretrained via discrete diffusion, while still using the same masked-modeling objective for customization as for the other methods.

Most notably, ProteinTTT enables accurate structure prediction for targets that are poorly predicted with the original models. For instance, Figure 1 presents a strongly improved structure predicted using ESMFold + ProteinTTT for the target that was part of the CASP14 competition and shown as an unsuccessful case in the original ESMFold publication (Lin et al. (2023), Fig. 2E). Another example is shown in Figure 3, where ProteinTTT refined the structure prediction from a low-quality prediction (TM-score = 0.29) to a nearly perfectly folded protein (TM-score = 0.92). Figure A4 shows that ESMFold + ProteinTTT maintains computational efficiency of ESMFold, being an order of magnitude faster than AlphaFold2. Figure A11 additionally demonstrates the robustness of ESM3 + ProteinTTT to the choice of hyperparameters.

4.2 PROTEIN FITNESS PREDICTION

The task of protein fitness prediction is to accurately order mutations of a protein based on their disruptive/favorable effects on protein functioning.

Evaluation Setup. We evaluate the models using ProteinGym, the state-of-the-art fitness prediction benchmark (Notin et al., 2024), focusing on its well-established zero-shot setup. Since the zero-shot

Table 2: **Customization with ProteinTTT improves protein fitness prediction.** The right section of the table presents performance averaged across individual proteins and then across different protein phenotypes, as classified in the ProteinGym benchmark (Notin et al., 2024). The middle column shows the final performance, averaged across all five phenotype classes. In total, ProteinGym contains 2.5 million mutations across 217 proteins. Standard deviations are calculated over 5 random seeds.

	Avg. Spearman \uparrow	Spearman by phenotype \uparrow				
		Activity	Binding	Expression	Organismal Fitness	Stability
ESM2 (35M) (Lin et al., 2023)	0.3211	0.3137	0.2907	0.3435	0.2184	0.4392
ESM2 (35M) + ProteinTTT (Ours)	0.3407 \pm 0.00014	0.3407	0.2942	0.3550	0.2403	0.4733
ProGen2-small (151M) (Nijkamp et al., 2023)	0.3255	0.3316	0.2681	0.3730	0.3283	0.3264
ProGen2-small (151M) + ProteinTTT (Ours)	0.3591 \pm 0.00021	0.3827	0.2960	0.3875	0.3302	0.3992
SaProt (35M) (Su et al., 2023)	0.4062	0.3721	0.3568	0.4390	0.2879	0.5749
SaProt (35M) + ProteinTTT (Ours)	0.4106 \pm 0.00004	0.3783	0.3569	0.4430	0.2955	0.5795
ESM2 (650M) (Lin et al., 2023)	0.4139	0.4254	0.3366	0.4151	0.3691	0.5233
ESM2 (650M) + ProteinTTT (Ours)	0.4153 \pm 0.00003	0.4323	0.3376	0.4168	0.3702	0.5195
SaProt (650M) (Su et al., 2023)	0.4569	0.4584	0.3785	0.4884	0.3670	0.5919
SaProt (650M) + ProteinTTT (Ours)	0.4583 \pm 0.00001	0.4593	0.3790	0.4883	0.3754	0.5896
ProSST (K=2048) (Li et al., 2024)	0.5068	0.4758	0.4448	0.5302	0.4306	0.6526
ProSST (K=2048) + ProteinTTT (Ours)	0.5087 \pm 0.00004	0.4822	0.4470	0.5321	0.4315	0.6507

setup only provides a test set without any data split, we also validate ProteinTTT on independent data. To achieve this, we create a new fitness prediction dataset mined from MaveDB, a public repository of Multiplexed Assays of Variant Effect (MAVEs) (Esposito et al., 2019). Following ProteinGym, we report Spearman correlation between predicted and experimental fitness values.

Results. ProteinTTT consistently enhances fitness prediction performance of all the tested models across varying model scales (35M and 650M parameters for both ESM2 and SaProt; 110M for ProSST) and both datasets, i.e., test ProteinGym (Table 2) and validation MaveDB (Table A5). Notably, ProSST + ProteinTTT sets a new state of the art on the ProteinGym benchmark (Spearman correlation coefficients calculated for individual deep mutational scanning experiments (DMSs) have statistically significant difference according to a paired t-test with $p < 0.05$).

We observe that ProteinTTT primarily improves performance for proteins with low MSA depth (i.e., the number of available homologous sequences), suggesting that single-sequence customization enhances predictions for proteins with fewer similar sequences in the training data (Table A4). The fact that ProteinTTT more effectively improves the performance of smaller ESM2 and SaProt models compared to their larger variants may be a result of the benchmark performance being saturated for larger models, consistent with a recent observation (Notin, 2025). We provide a qualitative example showing how ESM2 (650M) + ProteinTTT significantly improves fitness prediction by capturing residues critical for protein stability (Figure A5). We also demonstrate that customization can be combined with evolutionary information from MSA to further boost fitness prediction (Appendix C).

4.3 PROTEIN FUNCTION PREDICTION

Finally, we demonstrate a proof of concept for customization in the context of protein function prediction. We experiment with two tasks: predicting protein location within a cell (Stärk et al., 2021), and substrate classification for terpene synthases (TPS), enzymes producing the largest class of natural products (Samusevich et al., 2025). Appendix D shows that per-protein customization with ProteinTTT consistently enhances the performance of representative models on both tasks.

5 CASE STUDIES

ProteinTTT can be incorporated into structure, fitness, or function prediction pipelines with a few lines of code (Appendix E). Here, we demonstrate two challenging case studies: improving modeling of antibody–antigen loops (Section 5.1) and expanding known structures of viral proteins (Section 5.2).

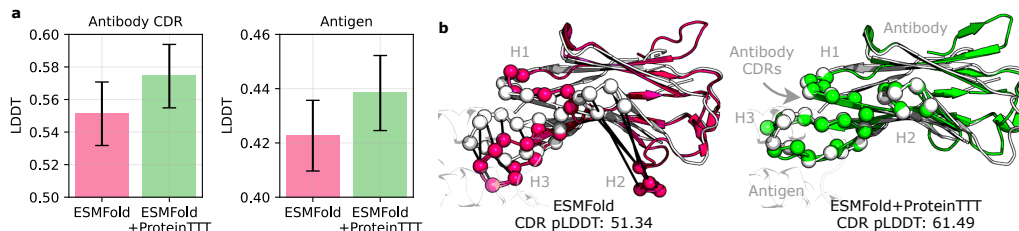


Figure 4: **ProteinTTT improves modeling of antibody-antigen loops.** (a) Average LDDT on the antibody complementarity-determining regions (CDRs, 175 structures) and antigens (814 structures) from the SABDab dataset with ESMFold pLDDT < 70. Error bars indicate 95% confidence intervals estimated from 1000 bootstrap samples. (b) Example of improved structure prediction for CDRs in the 8K2W entry. The CDR regions H1, H2, and H3, i.e., the parts of the antibody that bind to the antigen, are highlighted with spheres, while black lines show the alignment error between the ground-truth CDR structure (white) and the predictions (colored).

5.1 MODELING ANTIBODY-ANTIGEN LOOPS

Accurately predicting structures of antibodies (e.g., human defensive proteins) and antigens (e.g., viral proteins) enables rational design of new therapeutics (Bennett et al., 2025). However, the presence of highly variable loop regions makes modeling of these interactions a long-standing challenge. Here, we show that ProteinTTT substantially improves structure prediction for these loop-formed complementarity-determining regions (CDRs) of antibodies, i.e., the parts that bind antigens, as well as for antigens themselves, on the well-established SABDab dataset (Dunbar et al., 2014).

We take the structures from SABDab that are not predicted well by ESMFold (pLDDT < 70) and show that ProteinTTT improves the LDDT score for 115 of 175 antibody CDR substructures (66%) and 487 of 814 antigen chains (60%). As shown in Figure 4a, ESMFold + ProteinTTT achieves significantly higher average LDDT scores compared to general-purpose ESMFold (paired t-test p-value < 0.05). Figure 4b illustrates how ProteinTTT enables accurate prediction of all three CDRs in an antibody chain, providing an improved understanding of its binding interface with the corresponding antigen.

5.2 EXPANDING KNOWN STRUCTURES OF VIRAL PROTEINS

Predicting the structures of viral proteins is vital for vaccine development, antiviral design, and understanding infection (Bravi, 2024). Nevertheless, it remains challenging due to the high mutation rate, which often leaves viral proteins without close homologs or experimental structures in databases (Kim et al., 2025). Here, we demonstrate that per-protein customized predictions with ESMFold + ProteinTTT improve viral protein structure prediction, substantially expanding the Big Fantastic Virus Database—the comprehensive repository of 351,242 viral protein structures (Kim et al., 2025).

Among all the entries in BFVD, predicted with AlphaFold2 through ColabFold (Mirdita et al., 2022) using MSAs constructed from Logan (Chikhi et al., 2024), only 55% have high-quality structure predictions (pLDDT > 70). We apply ESMFold and ESMFold + ProteinTTT to the BFVD entries to expand the database with higher-quality structures. This is achieved by applying all three methods to the specific protein and taking the predicted structure with the highest pLDDT. While ESMFold manages to improve the predicted structure (as measured by pLDDT) for 10% of the BFVD proteins, ESMFold + ProteinTTT leads to an improvement for 19% of the dataset entries, substantially increasing the quality of known viral protein structures (Figure 5a).

We validate that the improved pLDDT confidence values from ESMFold + ProteinTTT correlate with the quality of the predicted structures, as measured by LDDT against reference AlphaFold2 structures having pLDDT > 90 (Pearson = 0.875; Figure A9). Notably, the largest improvements in pLDDT align with the largest improvements in LDDT (Figure 5b). We find that the benefit of customization saturates with the number of homologs available for a protein, indicating that ProteinTTT is most effective for challenging, out-of-distribution proteins (Figure 5c). Finally, Figure 5d–g shows examples where ProteinTTT enables high-confidence structure predictions in cases where general-purpose, uncustomized AlphaFold2 and ESMFold struggle.

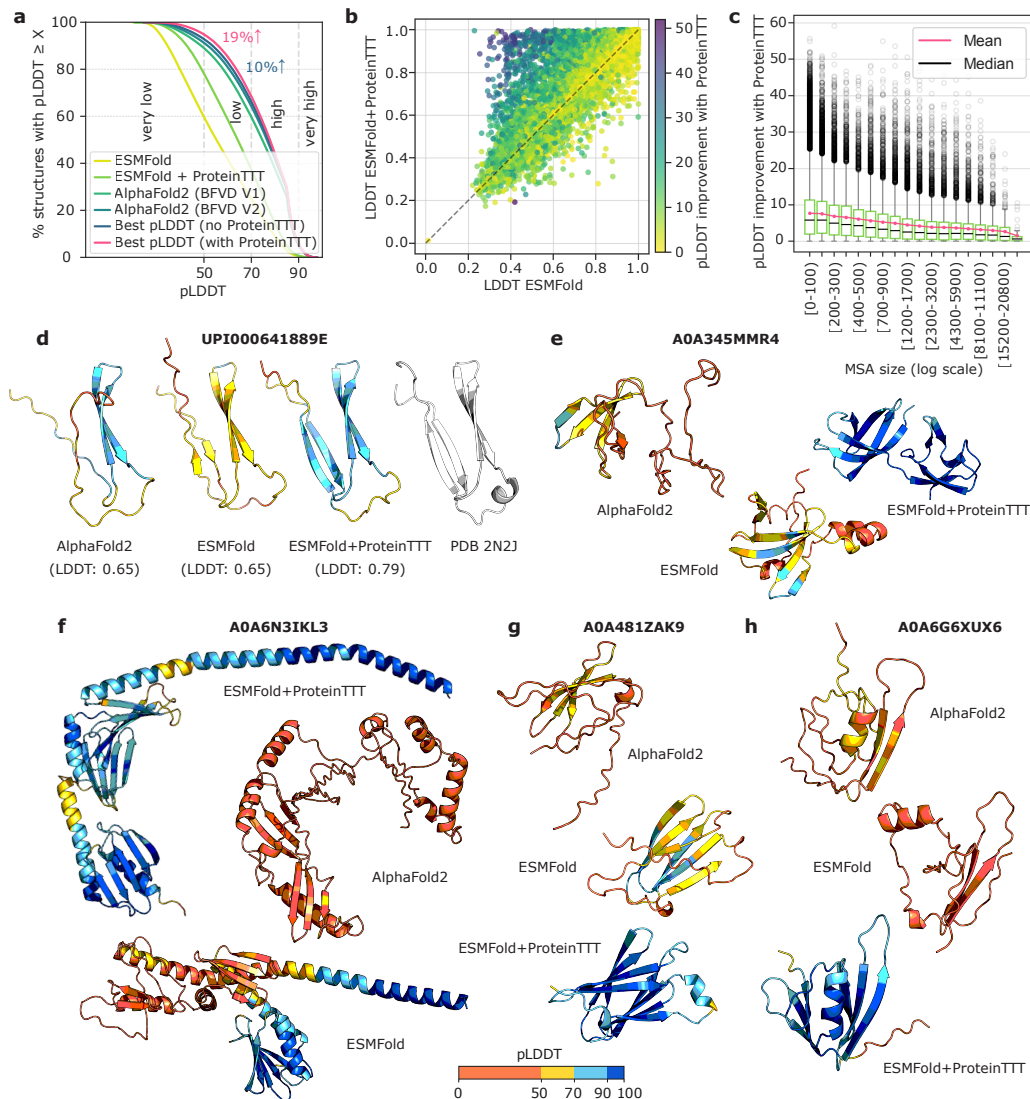


Figure 5: **ProteinTTT expands the Big Fantastic Virus Database (BFVD).** (a) ProteinTTT (light green) substantially improves the performance of ESMFold (yellow) on viral proteins, yielding better structures (pink) for 19% of BFVD entries compared to the original predictions by AlphaFold2 (green). (b) Improvements in pLDDT for ESMFold after ProteinTTT correspond to improvements in LDDT, as benchmarked against BFVD AlphaFold2 structures with pLDDT > 90 . (c) ProteinTTT provides the largest pLDDT improvements (y-axis) for the most out-of-distribution proteins, i.e., those with the smallest MSAs (left on the x-axis) from the Logan database. (d) Structural comparison for BFVD entry UPI000641889E against the PDB structure 2N2J (100% sequence identity) shows that ESMFold + ProteinTTT yields a prediction closest to the ground truth (gray), as also measured by LDDT. (e–g) Additional examples of high-quality viral structures (as measured by pLDDT) predicted with ESMFold + ProteinTTT but not with ESMFold or AlphaFold2. Higher pLDDT values are better.

6 DISCUSSION

We introduce ProteinTTT, a method for customizing protein language models to individual targets. ProteinTTT consistently improves performance across various models, their scales, and downstream tasks. It excels on challenging, out-of-distribution examples where general models often fail. We demonstrate its practical value through two case studies: enhancing the structural prediction of difficult antibody-antigen loops and improving 19% of low-confidence viral protein structures in the Big Fantastic Virus Database. Our work establishes per-protein customization as a powerful and practical tool for biological research.

REFERENCES

- 486
487
488 Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf
489 Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure
490 prediction of biomolecular interactions with alphafold 3. *Nature*, pp. 1–3, 2024. doi: 10.1038/
491 s41586-024-07487-w. URL <https://doi.org/10.1038/s41586-024-07487-w>.
- 492 Sarah Alamdari, Nitya Thakkar, Rianne Van Den Berg, Neil Tenenholtz, Robert Strome, Alan M
493 Moses, Alex X Lu, Nicolò Fusi, Ava P Amini, and Kevin K Yang. Protein generation with
494 evolutionary diffusion: sequence is all you need. *BioRxiv*, pp. 2023–09, 2023. doi: 10.1101/2023.
495 09.11.556673. URL <https://doi.org/10.1101/2023.09.11.556673>.
- 496 Ethan C Alley, Grigory Khimulya, Surojit Biswas, Mohammed AlQuraishi, and George M Church.
497 Unified rational protein engineering with sequence-based deep representation learning. *Nature*
498 *methods*, 16(12):1315–1322, 2019. doi: 10.1038/s41592-019-0598-1. URL <https://doi.org/10.1038/s41592-019-0598-1>.
- 500 José Juan Almagro Armenteros, Casper Kaae Sønderby, Søren Kaae Sønderby, Henrik Nielsen,
501 and Ole Winther. Deeploc: prediction of protein subcellular localization using deep learning.
502 *Bioinformatics*, 33(21):3387–3395, 2017. doi: 10.1093/bioinformatics/btx431. URL <https://doi.org/10.1093/bioinformatics/btx431>.
- 503 Frances M. Ashcroft, Matthew Lloyd, and Elizabeth A. Haythorne. Glucokinase activity in diabetes:
504 too much of a good thing? *Trends in Endocrinology & Metabolism*, 34(2):119–130, Feb 2023.
505 ISSN 1043-2760. doi: 10.1016/j.tem.2022.12.007. URL <https://doi.org/10.1016/j.tem.2022.12.007>.
- 506 Pratyay Banerjee, Tejas Gokhale, and Chitta Baral. Self-supervised test-time learning for reading
507 comprehension. *arXiv preprint arXiv:2103.11263*, 2021. doi: 10.48550/arXiv.2103.11263. URL
508 <https://doi.org/10.48550/arXiv.2103.11263>.
- 509 Eyal Ben-David, Nadav Oved, and Roi Reichart. Pada: Example-based prompt learning for on-the-fly
510 adaptation to unseen domains. *Transactions of the Association for Computational Linguistics*,
511 10:414–433, 2022. doi: 10.48550/arXiv.2102.12206. URL <https://doi.org/10.48550/arXiv.2102.12206>.
- 512 Nathaniel R Bennett, Joseph L Watson, Robert J Ragotte, Andrew J Borst, DéJenaé L See, Connor
513 Weidle, Riti Biswas, Yutong Yu, Ellen L Shrock, Russell Ault, et al. Atomically accurate de novo
514 design of antibodies with rfdiffusion. *bioRxiv*, pp. 2024–03, 2025. doi: 10.1101/2024.03.14.585103.
515 URL <https://doi.org/10.1101/2024.03.14.585103>.
- 516 Barbara Bravi. Development and use of machine learning algorithms in vaccine target selection. *npj*
517 *Vaccines*, 9(1):15, 2024. doi: 10.1038/s41541-023-00795-8. URL <https://doi.org/10.1038/s41541-023-00795-8>.
- 518 Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. doi:
519 10.48550/arXiv.2005.1416. URL <https://doi.org/10.48550/arXiv.2005.1416>.
- 520 Anton Bushuiev, Roman Bushuiev, Anatolii Filkin, Petr Kouba, Marketa Gabrielova, Michal Gabriel,
521 Jiri Sedlar, Tomas Pluskal, Jiri Damborsky, Stanislav Mazurenko, et al. Learning to design protein-
522 protein interactions with enhanced generalization. *arXiv preprint arXiv:2310.18515*, 2023. doi:
523 10.48550/arXiv.2310.18515. URL <https://arxiv.org/abs/2310.18515>.
- 524 Andrew Campbell, Jason Yim, Regina Barzilay, Tom Rainforth, and Tommi Jaakkola. Generative
525 flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design.
526 *arXiv preprint arXiv:2402.04997*, 2024. doi: 10.48550/arXiv.2402.04997. URL <https://doi.org/10.48550/arXiv.2402.04997>.
- 527 Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and
528 Tony Robinson. One billion word benchmark for measuring progress in statistical language
529 modeling. *arXiv preprint arXiv:1312.3005*, 2013. doi: 10.48550/arXiv.1312.3005. URL <https://doi.org/10.48550/arXiv.1312.3005>.
- 530
531
532
533
534
535
536
537
538
539

- 540 Tianlong Chen and Chengyue Gong. Hotprotein: A novel framework for protein thermostability
541 prediction and editing. *NeurIPS 2022*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=YDJRFWBMNby)
542 [id=YDJRFWBMNby](https://openreview.net/forum?id=YDJRFWBMNby).
543
- 544 Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Taylor Applebaum, Alexan-
545 der Pritzel, Lai Hong Wong, Michal Zielinski, Tobias Sargeant, et al. Accurate proteome-wide
546 missense variant effect prediction with alphamissense. *Science*, 381(6664):eadg7492, 2023. doi:
547 10.1126/science.adg7492. URL [https://www.science.org/doi/10.1126/science.](https://www.science.org/doi/10.1126/science.adg7492)
548 [adg7492](https://www.science.org/doi/10.1126/science.adg7492).
- 549 Zhixiang Chi, Li Gu, Tao Zhong, Huan Liu, Yuanhao Yu, Konstantinos N Plataniotis, and Yang
550 Wang. Adapting to distribution shift by visual domain prompt generation. *arXiv preprint*
551 *arXiv:2405.02797*, 2024. doi: 10.48550/arXiv.2405.02797. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2405.02797)
552 [48550/arXiv.2405.02797](https://doi.org/10.48550/arXiv.2405.02797).
553
- 554 Rayan Chikhi, Téo Lemane, Raphaël Loll-Krippleber, Mercè Montoliu-Nerin, Brice Raffestin,
555 Antonio Pedro Camargo, Carson J Miller, Mateus Bernabe Fiamenghi, Daniel Paiva Agostinho,
556 Sina Majidian, et al. Logan: planetary-scale genome assembly surveys life’s diversity. *bioRxiv*,
557 pp. 2024–07, 2024. doi: 10.1101/2024.07.30.605881. URL [https://doi.org/10.1101/](https://doi.org/10.1101/2024.07.30.605881)
558 [2024.07.30.605881](https://doi.org/10.1101/2024.07.30.605881).
- 559 Yehlin Cho, Martin Pacesa, Zhidian Zhang, Bruno E Correia, and Sergey Ovchinnikov. Boltzdesign1:
560 Inverting all-atom structure prediction model for generalized biomolecular binder design. *bioRxiv*,
561 pp. 2025–04, 2025. doi: 10.1101/2025.04.06.647261. URL [https://doi.org/10.1101/](https://doi.org/10.1101/2025.04.06.647261)
562 [2025.04.06.647261](https://doi.org/10.1101/2025.04.06.647261).
563
- 564 Cyrus Chothia and Arthur M Lesk. Canonical structures for the hypervariable regions of immunoglob-
565 ulins. *Journal of molecular biology*, 196(4):901–917, 1987. doi: 10.1016/0022-2836(87)90412-8.
566 URL [https://doi.org/10.1016/0022-2836\(87\)90412-8](https://doi.org/10.1016/0022-2836(87)90412-8).
- 567 David W. Christianson. Structural and chemical biology of terpenoid cyclases. *Chemical Reviews*,
568 117(17):11570–11648, Sep 2017. ISSN 0009-2665. doi: 10.1021/acs.chemrev.7b00287. URL
569 <https://doi.org/10.1021/acs.chemrev.7b00287>.
570
- 571 The UniProt Consortium. Uniprot: the universal protein knowledgebase in 2023. *Nucleic acids*
572 *research*, 51(D1):D523–D531, 2023. doi: 10.1093/nar/gkac1052. URL [https://doi.org/](https://doi.org/10.1093/nar/gkac1052)
573 [10.1093/nar/gkac1052](https://doi.org/10.1093/nar/gkac1052).
- 574 Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F Milles,
575 Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. Robust deep learning–based
576 protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022. doi: 10.1126/
577 [science.add2187](https://doi.org/10.1126/science.add2187). URL <https://doi.org/10.1126/science.add2187>.
578
- 579 Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv*
580 *preprint arXiv:1810.04805*, 2018. doi: 10.48550/arXiv.1810.04805. URL [https://doi.org/](https://doi.org/10.48550/arXiv.1810.04805)
581 [10.48550/arXiv.1810.04805](https://doi.org/10.48550/arXiv.1810.04805).
- 582 Daniel J Diaz, Chengyue Gong, Jeffrey Ouyang-Zhang, James M Loy, Jordan Wells, David Yang,
583 Andrew D Ellington, Alex Dimakis, and Adam R Klivans. Stability oracle: a structure-based
584 graph-transformer for identifying stabilizing mutations. *BioRxiv*, pp. 2023–05, 2023. doi: 10.1038/
585 [s41467-024-49780-2](https://doi.org/10.1038/s41467-024-49780-2). URL <https://doi.org/10.1038/s41467-024-49780-2>.
586
- 587 Henry Dieckhaus, Michael Brocidiacono, Nicholas Z Randolph, and Brian Kuhlman. Transfer learn-
588 ing to leverage larger datasets for improved prediction of protein stability changes. *Proceedings of*
589 *the National Academy of Sciences*, 121(6):e2314853121, 2024. doi: 10.1073/pnas.2314853121.
590 URL <https://doi.org/10.1073/pnas.2314853121>.
- 591 James Dunbar, Konrad Krawczyk, Jinwoo Leem, Terry Baker, Angelika Fuchs, Guy Georges, Jiye
592 Shi, and Charlotte M Deane. Sabdab: the structural antibody database. *Nucleic acids research*, 42
593 (D1):D1140–D1146, 2014. doi: 10.1093/nar/gkt1043. URL [https://doi.org/10.1093/](https://doi.org/10.1093/nar/gkt1043)
[nar/gkt1043](https://doi.org/10.1093/nar/gkt1043).

- 594 Janani Durairaj, Andrew M Waterhouse, Toomas Mets, Tetiana Brodiazhenko, Minhal Abdul-
595 lah, Gabriel Studer, Gerardo Tauriello, Mehmet Akdel, Antonina Andreeva, Alex Bateman,
596 et al. Uncovering new families and folds in the natural protein universe. *Nature*, 622(7983):
597 646–653, 2023. doi: 10.1038/s41586-023-06622-3. URL [https://doi.org/10.1038/
598 s41586-023-06622-3](https://doi.org/10.1038/s41586-023-06622-3).
- 599
600 Oliver Dutton, Sandro Bottaro, Istvan Redl, Michele Invernizzi, Albert Chung, Carlo Fisicaro, Falk
601 Hoffmann, Stefano Ruschetta, Fabio Airoidi, Louie Henderson, et al. Improving inverse folding
602 models at protein stability prediction without additional training or data. *bioRxiv*, pp. 2024–06,
603 2024. doi: 10.1101/2024.06.15.599145. URL [https://doi.org/10.1101/2024.06.15.
604 599145](https://doi.org/10.1101/2024.06.15.599145).
- 605 Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Yu Wang, Llion Jones, Tom
606 Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, et al. Prottrans: Toward understanding
607 the language of life through self-supervised learning. *IEEE transactions on pattern analysis
608 and machine intelligence*, 44(10):7112–7127, 2021. doi: 10.1109/tpami.2021.3095381. URL
609 <https://doi.org/10.1109/tpami.2021.3095381>.
- 610
611 Ahmed Elnaggar, Hazem Essam, Wafaa Salah-Eldin, Walid Moustafa, Mohamed Elkerdawy, Char-
612 lotte Rochereau, and Burkhard Rost. Ankh: Optimized protein language model unlocks general-
613 purpose modelling. *arXiv preprint arXiv:2301.06568*, 2023. doi: 10.48550/arXiv.2301.06568.
614 URL <https://doi.org/10.48550/arXiv.2301.06568>.
- 615 Daniel Esposito, Jochen Weile, Jay Shendure, Lea M Starita, Anthony T Papenfuss, Frederick P Roth,
616 Douglas M Fowler, and Alan F Rubin. Mavedb: an open-source platform to distribute and interpret
617 data from multiplexed assays of variant effect. *Genome biology*, 20:1–11, 2019. doi: 10.1186/
618 s13059-019-1845-6. URL <https://doi.org/10.1186/s13059-019-1845-6>.
- 619
620 Richard Evans, Michael O’Neill, Alexander Pritzel, Natasha Antropova, Andrew Senior, Tim Green,
621 Augustin Židek, Russ Bates, Sam Blackwell, Jason Yim, et al. Protein complex prediction with
622 alphafold-multimer. *biorxiv*, pp. 2021–10, 2021. doi: doi.org/10.1101/2021.10.04.463034. URL
623 <https://doi.org/10.1101/2021.10.04.463034>.
- 624 Xiaomin Fang, Fan Wang, Lihang Liu, Jingzhou He, Dayong Lin, Yingfei Xiang, Kunrui Zhu,
625 Xiaonan Zhang, Hua Wu, Hui Li, et al. A method for multiple-sequence-alignment-free pro-
626 tein structure prediction using a protein language model. *Nature Machine Intelligence*, 5(10):
627 1087–1096, 2023. doi: 10.1038/s42256-023-00721-6. URL [https://doi.org/10.1038/
628 s42256-023-00721-6](https://doi.org/10.1038/s42256-023-00721-6).
- 629
630 Tao Feng, Ziqi Gao, Jiaxuan You, Chenyi Zi, Yan Zhou, Chen Zhang, and Jia Li. Deep reinforcement
631 learning for modelling protein complexes. *arXiv preprint arXiv:2405.02299*, 2024. doi: 10.48550/
632 arXiv.2405.02299. URL <https://doi.org/10.48550/arXiv.2405.02299>.
- 633 Yossi Gandelsman, Yu Sun, Xinlei Chen, and Alexei Efros. Test-time training with masked au-
634 toencoders. *Advances in Neural Information Processing Systems*, 35:29374–29385, 2022. doi:
635 10.48550/arXiv.2209.07522. URL <https://doi.org/10.48550/arXiv.2209.07522>.
- 636
637 Cade Gordon, Amy X Lu, and Pieter Abbeel. Protein language model fitness is a matter of preference.
638 *bioRxiv*, pp. 2024–10, 2024. doi: 10.1101/2024.10.03.616542. URL [https://doi.org/10.
639 1101/2024.10.03.616542](https://doi.org/10.1101/2024.10.03.616542).
- 640
641 Jan Gorodkin. Comparing two k-category assignments by a k-category correlation coefficient.
642 *Computational biology and chemistry*, 28(5-6):367–374, 2004. doi: 10.1016/j.compbiolchem.2004.
643 09.006. URL <https://doi.org/10.1016/j.compbiolchem.2004.09.006>.
- 644
645 Xin Gu, Patrick Jouandin, Pranav V. Lalgudi, Rich Binari, Max L. Valenstein, Michael A. Reid,
646 Annamarie E. Allen, Nolan Kamitaki, Jason W. Locasale, Norbert Perrimon, and David M.
647 Sabatini. Sestrin mediates detection of and adaptation to low-leucine diets in drosophila. *Nature*,
608(7921):209–216, Aug 2022. ISSN 1476-4687. doi: 10.1038/s41586-022-04960-2. URL
<https://doi.org/10.1038/s41586-022-04960-2>.

- 648 Muhammet F. Gulen, Natasha Samson, Alexander Keller, Marius Schwabenland, Chong Liu, Selene
649 Glück, Vivek V. Thacker, Lucie Favre, Bastien Mangeat, Lona J. Kroese, Paul Krimpenfort, Marco
650 Prinz, and Andrea Ablasser. cgas–sting drives ageing-related inflammation and neurodegeneration.
651 *Nature*, 620(7973):374–380, Aug 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06373-1.
652 URL <https://doi.org/10.1038/s41586-023-06373-1>.
- 653 Kathryn H. Gunn and Saskia B. Neher. Structure of dimeric lipoprotein lipase reveals a pore adjacent
654 to the active site. *Nature Communications*, 14(1):2569, May 2023. ISSN 2041-1723. doi: 10.1038/
655 s41467-023-38243-9. URL <https://doi.org/10.1038/s41467-023-38243-9>.
- 656 Sarah Gurev, Noor Youssef, Navami Jain, and Debora S Marks. Variant effect prediction with
657 reliability estimation across priority viruses. *bioRxiv*, pp. 2025–08, 2025. doi: 10.1101/2025.08.
658 04.668549. URL <https://doi.org/10.1101/2025.08.04.668549>.
- 659 Moritz Hardt and Yu Sun. Test-time training on nearest neighbors for large language models. *arXiv*
660 *preprint arXiv:2305.18466*, 2023. doi: 10.48550/arXiv.2305.18466. URL [https://doi.org/
661 10.48550/arXiv.2305.18466](https://doi.org/10.48550/arXiv.2305.18466).
- 662 Tomas Hayes, Roshan Rao, Halil Akin, Nicholas J Sofroniew, Deniz Oktay, Zeming Lin, Robert
663 Verkuil, Vincent Q Tran, Jonathan Deaton, Marius Wiggert, et al. Simulating 500 million years
664 of evolution with a language model. *bioRxiv*, pp. 2024–07, 2024. doi: 10.1126/science.ads0018.
665 URL <https://www.science.org/doi/10.1126/science.ads0018>.
- 666 Michael Heinzinger and Burkhard Rost. Teaching ai to speak protein. *Current opinion in structural*
667 *biology*, 91:102986, 2025. doi: 10.1016/j.sbi.2025.102986. URL [https://doi.org/10.
668 1016/j.sbi.2025.102986](https://doi.org/10.1016/j.sbi.2025.102986).
- 669 Lucas Torroba Hennigen and Yoon Kim. Deriving language models from masked language models.
670 *arXiv preprint arXiv:2305.15501*, 2023. doi: 10.48550/arXiv.2305.15501. URL [https://doi.
671 org/10.48550/arXiv.2305.15501](https://doi.org/10.48550/arXiv.2305.15501).
- 672 Thomas A Hopf, John B Ingraham, Frank J Poelwijk, Charlotta PI Schärfe, Michael Springer, Chris
673 Sander, and Debora S Marks. Mutation effects predicted from sequence co-variation. *Nature*
674 *biotechnology*, 35(2):128–135, 2017. doi: 10.1038/nbt.3769. URL [https://doi.org/10.
675 1038/nbt.3769](https://doi.org/10.1038/nbt.3769).
- 676 Gerta Hoxhaj and Brendan D. Manning. The pi3k–akt network at the interface of oncogenic signalling
677 and cancer metabolism. *Nature Reviews Cancer*, 20(2):74–88, Feb 2020. ISSN 1474-1768. doi: 10.
678 1038/s41568-019-0216-7. URL <https://doi.org/10.1038/s41568-019-0216-7>.
- 679 Cheng-Yen Hsieh, Xinyou Wang, Daiheng Zhang, Dongyu Xue, Fei Ye, Shujian Huang, Zaixiang
680 Zheng, and Quanquan Gu. Elucidating the design space of multimodal protein language models.
681 *arXiv preprint arXiv:2504.11454*, 2025. doi: 10.48550/arXiv.2504.11454. URL [https://doi.
682 org/10.48550/arXiv.2504.11454](https://doi.org/10.48550/arXiv.2504.11454).
- 683 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
684 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
685 *arXiv:2106.09685*, 2021. doi: 10.48550/arXiv.2106.09685. URL [https://doi.org/10.
686 48550/arXiv.2106.09685](https://doi.org/10.48550/arXiv.2106.09685).
- 687 Jonas Hübötter, Sascha Bongni, Ido Hakimi, and Andreas Krause. Efficiently learning at test-time:
688 Active fine-tuning of llms. *arXiv preprint arXiv:2410.08020*, 2024. doi: 10.48550/arXiv.2410.
689 08020. URL <https://doi.org/10.48550/arXiv.2410.08020>.
- 690 Jonas Hübötter, Patrik Wolf, Alexander Shevchenko, Dennis Jüni, Andreas Krause, and Gil Kur.
691 Specialization after generalization: Towards understanding test-time training in foundation models.
692 *arXiv preprint arXiv:2509.24510*, 2025. doi: 10.48550/arXiv.2509.24510. URL [https://doi.
693 org/10.48550/arXiv.2509.24510](https://doi.org/10.48550/arXiv.2509.24510).
- 694 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,
695 Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al. Highly accurate
696 protein structure prediction with alphafold. *nature*, 596(7873):583–589, 2021. doi: 10.1038/
697 s41586-021-03819-2. URL <https://doi.org/10.1038/s41586-021-03819-2>.
- 698
699
700
701

- 702 Pranav Kantroo, Gunter Wagner, and Benjamin Machta. Pseudo-perplexity in one fell swoop for
703 protein fitness estimation. *bioRxiv*, pp. 2024–07, 2024. doi: 10.48550/arXiv.2407.07265. URL
704 <https://doi.org/10.48550/arXiv.2407.07265>.
705
- 706 Neerav Karani, Ertunc Erdil, Krishna Chaitanya, and Ender Konukoglu. Test-time adaptable neural
707 networks for robust medical image segmentation. *Medical Image Analysis*, 68:101907, 2021.
708 doi: 10.1016/j.media.2020.101907. URL [https://doi.org/10.1016/j.media.2020.](https://doi.org/10.1016/j.media.2020.101907)
709 101907.
- 710 Zuzana Keckesova, Joana Liu Donaher, Jasmine De Cock, Elizaveta Freinkman, Susanne Lingrell,
711 Daniel A. Bachovchin, Brian Bierie, Verena Tischler, Aurelia Noske, Marian C. Okondo, Ferenc
712 Reinhardt, Prathapan Thiru, Todd R. Golub, Jean E. Vance, and Robert A. Weinberg. Lactb is a
713 tumour suppressor that modulates lipid metabolism and cell state. *Nature*, 543(7647):681–686,
714 Mar 2017. ISSN 1476-4687. doi: 10.1038/nature21408. URL [https://doi.org/10.1038/](https://doi.org/10.1038/nature21408)
715 nature21408.
- 716 Rachel Seongeun Kim, Eli Levy Karin, Milot Mirdita, Rayan Chikhi, and Martin Steinegger. Bfvd—a
717 large repository of predicted viral protein structures. *Nucleic Acids Research*, 53(D1):D340–D347,
718 2025. doi: 10.1093/nar/gkae1119. URL <https://doi.org/10.1093/nar/gkae1119>.
719
- 720 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International*
721 *Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015,*
722 *Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6980>.
723
- 724 Andrew Kirjner, Jason Yim, Raman Samusevich, Shahar Bracha, Tommi S Jaakkola, Regina Barzilay,
725 and Ila R Fiete. Improving protein optimization with smoothed fitness landscapes. In *The Twelfth*
726 *International Conference on Learning Representations*, 2023. URL [https://openreview.](https://openreview.net/forum?id=rx1F2Zv8x0)
727 net/forum?id=rx1F2Zv8x0.
- 728 Petr Kouba, Pavel Kohout, Faraneh Haddadi, Anton Bushuiev, Raman Samusevich, Jiri Sedlar, Jiri
729 Damborsky, Tomas Pluskal, Josef Sivic, and Stanislav Mazurenko. Machine learning-guided
730 protein engineering. *ACS catalysis*, 13(21):13863–13895, 2023. doi: 10.1021/acscatal.3c02743.
731 URL <https://doi.org/10.1021/acscatal.3c02743>.
- 732 Elodie Laine, Yasaman Karami, and Alessandra Carbone. Gemme: a simple and fast global epistatic
733 model predicting mutational effects. *Molecular biology and evolution*, 36(11):2604–2619, 2019.
734 doi: 10.1093/molbev/msz179. URL <https://doi.org/10.1093/molbev/msz179>.
735
- 736 Mingchen Li, Yang Tan, Xinzhu Ma, Bozitao Zhong, Huiqun Yu, Ziyi Zhou, Wanli
737 Ouyang, Bingxin Zhou, Liang Hong, and Pan Tan. Prosst: Protein language model-
738 eling with quantized structure and disentangled attention. *bioRxiv*, pp. 2024–04,
739 2024. URL [https://proceedings.neurips.cc/paper_files/paper/2024/](https://proceedings.neurips.cc/paper_files/paper/2024/file/3ed57b293db0aab7cc30c44f45262348-Paper-Conference.pdf)
740 file/3ed57b293db0aab7cc30c44f45262348-Paper-Conference.pdf.
- 741 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,
742 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom
743 Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level pro-
744 tein structure with a language model. *Science*, 379(6637):1123–1130, 2023. doi: 10.1126/
745 science.ade2574. URL [https://www.science.org/doi/abs/10.1126/science.](https://www.science.org/doi/abs/10.1126/science.ade2574)
746 ade2574.
- 747 Yuejiang Liu, Parth Kothari, Bastien van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and
748 Alexandre Alahi. TTT++: when does self-supervised test-time training fail or thrive? In
749 *Advances in Neural Information Processing Systems 34: Annual Conference on Neural*
750 *Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp.
751 21808–21820, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/](https://proceedings.neurips.cc/paper/2021/hash/b618c3210e934362ac261db280128c22-Abstract.html)
752 b618c3210e934362ac261db280128c22-Abstract.html.
753
- 754 Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):
755 129–137, 1982. doi: 10.1109/TIT.1982.1056489. URL [https://doi.org/10.1109/TIT.](https://doi.org/10.1109/TIT.1982.1056489)
1982.1056489.

- 756 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International*
757 *Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*.
758 OpenReview.net, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
759
- 760 Valerio Mariani, Marco Biasini, Alessandro Barbato, and Torsten Schwede. lddt: a local
761 superposition-free score for comparing protein structures and models using distance difference
762 tests. *Bioinformatics*, 29(21):2722–2728, 2013. doi: 10.1093/bioinformatics/btt473. URL
763 <https://doi.org/10.1093/bioinformatics/btt473>.
764
- 765 Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, and Alex Rives.
766 Language models enable zero-shot prediction of the effects of mutations on pro-
767 tein function. *Advances in neural information processing systems*, 34:29287–29303,
768 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/
f51338d736f95dd42427296047067694-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/f51338d736f95dd42427296047067694-Abstract.html).
769
- 770 Peter Mikhael, Itamar Chinn, and Regina Barzilay. Clipzyme: Reaction-conditioned virtual screening
771 of enzymes. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Aus-*
772 *tria, July 21-27, 2024*. OpenReview.net, 2024. URL [https://openreview.net/forum?
id=0mYAK6Yhhm](https://openreview.net/forum?id=0mYAK6Yhhm).
773
- 774 Milot Mirdita, Konstantin Schütze, Yoshitaka Moriwaki, Lim Heo, Sergey Ovchinnikov, and Mar-
775 tin Steinegger. Colabfold: making protein folding accessible to all. *Nature methods*, 19(6):
776 679–682, 2022. doi: 10.1038/s41592-022-01488-1. URL [https://doi.org/10.1038/
s41592-022-01488-1](https://doi.org/10.1038/s41592-022-01488-1).
777
- 778 Erik Nijkamp, Jeffrey A. Ruffolo, Eli N. Weinstein, Nikhil Naik, and Ali Madani. Progen2: Exploring
779 the boundaries of protein language models. *Cell Systems*, 14(11):968–978.e3, Nov 2023. ISSN
780 2405-4712. doi: 10.1016/j.cels.2023.10.002. URL [https://doi.org/10.1016/j.cels.
2023.10.002](https://doi.org/10.1016/j.cels.2023.10.002).
781
- 782 Pascal Notin. Have we hit the scaling wall for protein language models? Sub-
783 stack blog post, May 7 2025. URL [https://pascalnotin.substack.com/p/
have-we-hit-the-scaling-wall-for](https://pascalnotin.substack.com/p/have-we-hit-the-scaling-wall-for).
784
- 785 Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood Van Niekerk, Steffanie Paul, Han Spinner, Nathan
786 Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, et al. Proteingym: Large-scale benchmarks
787 for protein fitness prediction and design. *Advances in Neural Information Processing Systems*,
788 36, 2024. doi: 10.1101/2023.12.07.570727. URL [https://doi.org/10.1101/2023.12.
07.570727](https://doi.org/10.1101/2023.12.07.570727).
789
- 790 Dong oh Seo, David O’Donnell, Nimansha Jain, Jason D. Ulrich, Jasmin Herz, Yuhao Li, Mackenzie
791 Lemieux, Jiye Cheng, Hao Hu, Javier R. Serrano, Xin Bao, Emily Franke, Maria Karlsson,
792 Martin Meier, Su Deng, Chandani Desai, Hemraj Dodiya, Janaki Lelwala-Guruge, Scott A.
793 Handley, Jonathan Kipnis, Sangram S. Sisodia, Jeffrey I. Gordon, and David M. Holtzman.
794 Apoe isoform- and microbiota-dependent progression of neurodegeneration in a mouse model of
795 tauopathy. *Science*, 379(6628):eadd1236, 2023. doi: 10.1126/science.add1236. URL <https://www.science.org/doi/abs/10.1126/science.add1236>.
796
- 797 Andrei Papkou, Lucia Garcia-Pastor, José Antonio Escudero, and Andreas Wagner. A rugged yet eas-
798 ily navigable fitness landscape. *Science*, 382(6673):eadh3860, 2023. doi: 10.1126/science.adh3860.
799 URL <https://www.science.org/doi/abs/10.1126/science.adh3860>.
800
- 801 Saro Passaro, Gabriele Corso, Jeremy Wohlwend, Mateo Reveiz, Stephan Thaler, Vignesh Ram
802 Somnath, Noah Getz, Tally Portnoi, Julien Roy, Hannes Stark, et al. Boltz-2: Towards accurate
803 and efficient binding affinity prediction. *BioRxiv*, 2025. doi: 10.1101/2025.06.14.659707. URL
804 <https://doi.org/10.1101/2025.06.14.659707>.
805
- 806 A Paszke. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint*
807 *arXiv:1912.01703*, 2019. doi: 10.48550/arXiv.1912.01703. URL [https://doi.org/10.
48550/arXiv.1912.01703](https://doi.org/10.48550/arXiv.1912.01703).
808
- 809

- 810 Predrag Radivojac and et al. A large-scale evaluation of computational protein function prediction.
811 *Nature Methods*, 10(3):221–227, Mar 2013. ISSN 1548-7105. doi: 10.1038/nmeth.2340. URL
812 <https://doi.org/10.1038/nmeth.2340>.
813
- 814 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark
815 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference*
816 *on machine learning*, pp. 8821–8831. Pmlr, 2021. doi: 10.48550/arXiv.2102.12092. URL
817 <https://doi.org/10.48550/arXiv.2102.12092>.
- 818 Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter
819 Abbeel, and Yun Song. Evaluating protein transfer learning with tape. *Advances in neural*
820 *information processing systems*, 32, 2019. doi: 10.48550/arXiv.1906.08230. URL <https://doi.org/10.48550/arXiv.1906.08230>.
821
- 822 Roshan Rao, Joshua Meier, Tom Sercu, Sergey Ovchinnikov, and Alexander Rives. Transformer
823 protein language models are unsupervised structure learners. *Biorxiv*, pp. 2020–12, 2020. doi:
824 10.1101/2020.12.15.422761. URL <https://doi.org/10.1101/2020.12.15.422761>.
825
- 826 Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu, and
827 Alexander Rives. Msa transformer. In *International Conference on Machine Learning*, pp. 8844–
828 8856. PMLR, 2021. URL <https://proceedings.mlr.press/v139/rao21a.html>.
829
- 830 Ali Razavi, Aäron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with
831 VQ-VAE-2. In *Advances in Neural Information Processing Systems 32: Annual Conference on*
832 *Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC,*
833 *Canada*, pp. 14837–14847, 2019. URL [https://proceedings.neurips.cc/paper/](https://proceedings.neurips.cc/paper/2019/hash/5f8e2fa1718d1bbcadf1cd9c7a54fb8c-Abstract.html)
834 [2019/hash/5f8e2fa1718d1bbcadf1cd9c7a54fb8c-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/5f8e2fa1718d1bbcadf1cd9c7a54fb8c-Abstract.html).
- 835 Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo,
836 Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from
837 scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National*
838 *Academy of Sciences*, 118(15):e2016239118, 2021. doi: 10.1073/pnas.2016239118. URL <https://doi.org/10.1073/pnas.2016239118>.
839
- 840 Xavier Robin, Juergen Haas, Rafal Gumienny, Anna Smolinski, Gerardo Tauriello, and Torsten
841 Schwede. Continuous automated model evaluation (cameo)—perspectives on the future of fully
842 automated evaluation of structure prediction methods. *Proteins: Structure, Function, and Bioin-*
843 *formatics*, 89(12):1977–1986, 2021. doi: 10.1002/prot.26213. URL [https://doi.org/10.](https://doi.org/10.1002/prot.26213)
844 [1002/prot.26213](https://doi.org/10.1002/prot.26213).
845
- 846 Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint*
847 *arXiv:1609.04747*, 2016. doi: 10.48550/arXiv.1609.04747. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.1609.04747)
848 [48550/arXiv.1609.04747](https://doi.org/10.48550/arXiv.1609.04747).
- 849 Julian Salazar, Davis Liang, Toan Q Nguyen, and Katrin Kirchhoff. Masked language model
850 scoring. *arXiv preprint arXiv:1910.14659*, 2019. doi: 10.18653/v1/2020.acl-main.240. URL
851 <https://doi.org/10.18653/v1/2020.acl-main.240>.
852
- 853 Raman Samusevich, Téó Hebra, Roman Bushuiev, Martin Engst, Jonáš Kulhánek, Anton Bushuiev,
854 Joshua D. Smith, Tereza Čalounová, Helena Smrčková, Marina Molineris, Renana Schwartz, Adéla
855 Tajovská, Milana Perković, Ratthachat Chatpatanasiri, Sotirios C. Kampranis, Dan Thomas Major,
856 Josef Sivic, and Tomáš Pluskal. Structure-enabled enzyme function prediction unveils elusive
857 terpenoid biosynthesis in archaea. *bioRxiv*, 2025. doi: 10.1101/2024.01.29.577750. URL <https://www.biorxiv.org/content/early/2025/04/29/2024.01.29.577750>.
858
- 859 Nicolae Sapoval, Amirali Aghazadeh, Michael G. Nute, Dinler A. Antunes, Advait Balaji, Richard
860 Baraniuk, C. J. Barberan, Ruth Dannenfelser, Chen Dun, Mohammadamin Edrisi, R. A. Leo
861 Elworth, Bryce Kille, Anastasios Kyrillidis, Luay Nakhleh, Cameron R. Wolfe, Zhi Yan, Vicky
862 Yao, and Todd J. Treangen. Current progress and open challenges for applying deep learning across
863 the biosciences. *Nature Communications*, 13(1):1728, Apr 2022. ISSN 2041-1723. doi: 10.1038/
s41467-022-29268-7. URL <https://doi.org/10.1038/s41467-022-29268-7>.

- 864 Robert Schmirler, Michael Heinzinger, and Burkhard Rost. Fine-tuning protein language models
865 boosts predictions across diverse tasks. *Nature Communications*, 15(1):7407, 2024. doi: 10.1038/
866 s41467-024-51844-2. URL <https://doi.org/10.1038/s41467-024-51844-2>.
867
- 868 Emre Sevgen, Joshua Moller, Adrian Lange, John Parker, Sean Quigley, Jeff Mayer, Poonam
869 Srivastava, Sitaram Gayatri, David Hosfield, Maria Korshunova, et al. Prot-vae: protein transformer
870 variational autoencoder for functional protein design. *bioRxiv*, pp. 2023–01, 2023. doi: 10.1073/
871 pnas.2408737122. URL <https://doi.org/10.1073/pnas.2408737122>.
- 872 Richard W Shuai, Talal Widatalla, Po-Ssu Huang, and Brian L Hie. Sidechain conditioning and
873 modeling for full-atom protein sequence design with fampnn. *bioRxiv*, pp. 2025–02, 2025. doi:
874 10.1101/2025.02.13.637498. URL <https://doi.org/10.1101/2025.02.13.637498>.
- 875 Elana Simon and James Zou. Interplm: Discovering interpretable features in protein language models
876 via sparse autoencoders. *Nature Methods*, pp. 1–11, 2025. doi: 10.1038/s41592-025-02836-7.
877 URL <https://doi.org/10.1038/s41592-025-02836-7>.
878
- 879 Peter Škrinjar, Jérôme Eberhardt, Janani Durairaj, and Torsten Schwede. Have protein-ligand co-
880 folding methods moved beyond memorisation? *BioRxiv*, pp. 2025–02, 2025. doi: 10.1101/2025.
881 02.03.636309. URL <https://doi.org/10.1101/2025.02.03.636309>.
- 882 Yidong Song, Qianmu Yuan, Sheng Chen, Yuansong Zeng, Huiying Zhao, and Yuedong Yang.
883 Accurately predicting enzyme functions through geometric graph learning on esmfold-predicted
884 structures. *Nature Communications*, 15(1):8180, 2024. doi: 10.1038/s41467-024-52533-w. URL
885 <https://doi.org/10.1038/s41467-024-52533-w>.
- 886 Hannes Stärk, Christian Dallago, Michael Heinzinger, and Burkhard Rost. Light attention predicts
887 protein location from the language of life. *Bioinformatics Advances*, 1(1):vbab035, 11 2021. ISSN
888 2635-0041. doi: 10.1093/bioadv/vbab035. URL [https://doi.org/10.1093/bioadv/
889 vbab035](https://doi.org/10.1093/bioadv/vbab035).
- 890 Jin Su, Chenchen Han, Yuyang Zhou, Junjie Shan, Xibin Zhou, and Fajie Yuan. Saprot: Protein
891 language modeling with structure-aware vocabulary. *bioRxiv*, pp. 2023–10, 2023. URL [https://
892 //openreview.net/forum?id=6MRm3G4NiU](https://openreview.net/forum?id=6MRm3G4NiU).
893
- 894 Sriram Subramaniam and Gerard J. Kleywegt. A paradigm shift in structural biology. *Nature*
895 *Methods*, 19(1):20–23, Jan 2022. ISSN 1548-7105. doi: 10.1038/s41592-021-01361-7. URL
896 <https://doi.org/10.1038/s41592-021-01361-7>.
- 897 Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training
898 with self-supervision for generalization under distribution shifts. In *International conference*
899 *on machine learning*, pp. 9229–9248. PMLR, 2020. URL [https://proceedings.mlr.
900 press/v119/sun20b/sun20b.pdf](https://proceedings.mlr.press/v119/sun20b/sun20b.pdf).
- 901 Nataša Tagasovska, Ji Won Park, Matthieu Kirchmeyer, Nathan C Frey, Andrew Martin Watkins,
902 Aya Abdelsalam Ismail, Arian Rokkum Jamasb, Edith Lee, Tyler Bryson, Stephen Ra, et al.
903 Antibody domainbed: Out-of-distribution generalization in therapeutic protein design. *arXiv*
904 *preprint arXiv:2407.21028*, 2024. doi: 10.48550/arXiv.2407.21028. URL [https://doi.org/
905 10.48550/arXiv.2407.21028](https://doi.org/10.48550/arXiv.2407.21028).
- 906 Kotaro Tsuboyama, Justas Dauparas, Jonathan Chen, Elodie Laine, Yasser Mohseni Behbahani,
907 Jonathan J Weinstein, Niall M Mangan, Sergey Ovchinnikov, and Gabriel J Rocklin. Mega-scale
908 experimental analysis of protein folding stability in biology and design. *Nature*, 620(7973):
909 434–444, 2023. doi: 10.1038/s41586-023-06328-6. URL [https://doi.org/10.1038/
911 s41586-023-06328-6](https://doi.org/10.1038/
910 s41586-023-06328-6).
- 912 Mike Tyers and Matthias Mann. From genomics to proteomics. *Nature*, 422(6928):193–197, Mar
913 2003. ISSN 1476-4687. doi: 10.1038/nature01510. URL [https://doi.org/10.1038/
915 nature01510](https://doi.org/10.1038/
914 nature01510).
- 916 Michel van Kempen, Stephanie S Kim, Charlotte Tumescheit, Milot Mirdita, Cameron LM Gilchrist,
917 Johannes Söding, and Martin Steinegger. Foldseek: fast and accurate protein structure search.
Biorxiv, pp. 2022–02, 2022. doi: 10.1038/s41587-023-01773-0. URL [https://doi.org/10.
1038/s41587-023-01773-0](https://doi.org/10.1038/s41587-023-01773-0).

- 918 Mihaly Varadi, Stephen Anyango, Mandar Deshpande, Sreenath Nair, Cindy Natassia, Galabina
919 Yordanova, David Yuan, Oana Stroe, Gemma Wood, Agata Laydon, et al. Alphafold protein
920 structure database: massively expanding the structural coverage of protein-sequence space with
921 high-accuracy models. *Nucleic acids research*, 50(D1):D439–D444, 2022. doi: 10.1093/nar/
922 gkab1061. URL <https://doi.org/10.1093/nar/gkab1061>.
- 923 A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. doi:
924 10.48550/arXiv.1706.03762. URL <https://doi.org/10.48550/arXiv.1706.03762>.
- 925
926 Renhao Wang, Yu Sun, Yossi Gandelsman, Xinlei Chen, Alexei A Efros, and Xiaolong Wang. Test-
927 time training on video streams. *arXiv preprint arXiv:2307.05014*, 2023. doi: 10.48550/arXiv.2307.
928 05014. URL <https://doi.org/10.48550/arXiv.2307.05014>.
- 929 Xinyou Wang, Zaixiang Zheng, Fei Ye, Dongyu Xue, Shujian Huang, and Quanquan Gu. Diffusion
930 language models are versatile protein learners. *arXiv preprint arXiv:2402.18567*, 2024a. doi:
931 10.48550/arXiv.2402.18567. URL <https://doi.org/10.48550/arXiv.2402.18567>.
- 932 Xinyou Wang, Zaixiang Zheng, Fei Ye, Dongyu Xue, Shujian Huang, and Quanquan Gu. Dplm-2:
933 A multimodal diffusion protein language model. *arXiv preprint arXiv:2410.13782*, 2024b. doi:
934 10.48550/arXiv.2410.13782. URL <https://doi.org/10.48550/arXiv.2410.13782>.
- 935
936 Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach,
937 Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein
938 structure and function with rfdiffusion. *Nature*, 620(7976):1089–1100, 2023. doi: 10.1038/
939 s41586-023-06415-8. URL <https://doi.org/10.1038/s41586-023-06415-8>.
- 940 Zehao Xiao, Xiantong Zhen, Ling Shao, and Cees GM Snoek. Learning to generalize across domains
941 on single test samples. *arXiv preprint arXiv:2202.08045*, 2022. doi: 10.48550/arXiv.2202.08045.
942 URL <https://doi.org/10.48550/arXiv.2202.08045>.
- 943 Tianhao Yu, Haiyang Cui, Jianan Canal Li, Yunan Luo, Guangde Jiang, and Huimin Zhao. En-
944 zyme function prediction using contrastive learning. *Science*, 379(6639):1358–1363, 2023.
945 doi: 10.1126/science.adf2465. URL [https://www.science.org/doi/abs/10.1126/
946 science.adf2465](https://www.science.org/doi/abs/10.1126/science.adf2465).
- 947
948 Yang Zhang and Jeffrey Skolnick. Scoring function for automated assessment of protein structure
949 template quality. *Proteins: Structure, Function, and Bioinformatics*, 57(4):702–710, 2004. doi:
950 10.1002/prot.20264. URL <https://doi.org/10.1002/prot.20264>.
- 951 Zhidian Zhang, Hannah K. Wayment-Steele, Garyk Brixli, Haobo Wang, Dorothee Kern, and Sergey
952 Ovchinnikov. Protein language models learn evolutionary statistics of interacting sequence
953 motifs. *Proceedings of the National Academy of Sciences*, 121(45):e2406285121, 2024. doi:
954 10.1073/pnas.2406285121. URL [https://www.pnas.org/doi/abs/10.1073/pnas.
955 2406285121](https://www.pnas.org/doi/abs/10.1073/pnas.2406285121).
- 956 Hao Zhao, Yuejiang Liu, Alexandre Alahi, and Tao Lin. On pitfalls of test-time adaptation. In
957 *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii,
958 USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 42058–42080. PMLR, 2023.
959 URL <https://proceedings.mlr.press/v202/zhao23d.html>.
- 960
961
962
963
964
965
966
967
968
969
970
971

972	APPENDIX	
973		
974	CONTENTS	
975		
976		
977	A Justification of customization via perplexity minimization	20
978		
979	B Customization beyond masked language modeling	20
980		
981	C Customization with multiple sequence alignment (MSA)	21
982		
983		
984	D Customization for protein function prediction	22
985		
986	E Implementation details	23
987		
988	F Experimental details	25
989	F.1 Protein structure prediction	25
990	F.1.1 Datasets	25
991	F.1.2 Metrics	26
992	F.1.3 Models	26
993	F.2 Protein fitness prediction	27
994	F.2.1 Datasets	27
995	F.2.2 Metrics	29
996	F.2.3 Models	29
997	F.3 Protein function prediction	31
998	F.3.1 Datasets	31
999	F.3.2 Metrics	31
1000	F.3.3 Models	32
1001		
1002		
1003		
1004		
1005		
1006		
1007	G Case study details	32
1008	G.1 Modeling antibody-antigen loops	32
1009	G.2 Expanding known structures of viral proteins	33
1010		
1011		
1012	H Extended results	33
1013	H.1 Detailed test performance	33
1014	H.2 Validation performance	34
1015	H.3 Runtime performance	34
1016		
1017		
1018		
1019	I Limitations and future work	34
1020		
1021		
1022		
1023		
1024		
1025		

A JUSTIFICATION OF CUSTOMIZATION VIA PERPLEXITY MINIMIZATION

While the paradigm of test-time customization has been investigated in other domains, the reasons behind its surprising effectiveness are not completely clear (Liu et al., 2021; Zhao et al., 2023). Here, we offer a potential justification for the effectiveness of ProteinTTT by linking it to perplexity minimization.

Perplexity has traditionally been used in natural language processing to evaluate how well models comprehend sentences (Brown, 2020; Chelba et al., 2013). Protein language modeling has adopted this metric to assess how effectively models “understand” amino acid sequences (Hayes et al., 2024; Lin et al., 2023). For bidirectional, random masking language models, which are the focus of this study, we consider the following definition of perplexity²:

$$\text{Perplexity}(x) = \exp\left(\frac{1}{|x|} \sum_{i=1}^{|x|} -\log p(x_i|x_{\setminus i}; \theta)\right), \quad (3)$$

where $|x|$ is the length of the input protein sequence x and $p(x_i|x_{\setminus i}; \theta)$ represents the probability that the model correctly predicts the token x_i at position i when it is masked on the input $x_{\setminus i}$. Perplexity ranges from 1 to infinity (the lower, the better), providing an intuitive measure of how well a model fits, on average, tokens in a given sequence. A perplexity value of 1 indicates that the model perfectly fits the sequence, accurately predicting all the true tokens.

Several studies have shown that lower perplexity on held-out protein sequences (calculated through the self-supervised track $g \circ f$) correlates with better performance on downstream tasks (via the supervised track $h \circ f$), such as predicting protein contacts (Rao et al., 2020), structure (Lin et al., 2023), or fitness (Kantroo et al., 2024). To give an example, we analyze the correlation between perplexity and structure prediction quality (Figure A1; see Section 4.1 for experimental details). A notable correlation suggests that reducing a model’s perplexity on a single target sample x (applied independently to all test samples) can lead to improved predictions on the downstream task (Figure 3; Figure A10).

Since we assume only a single target example x , the minimization of the masked language modeling loss $\mathcal{L}(x; \theta)$ (Equation (2)) on this example is directly linked to minimizing the perplexity $\text{Perplexity}(x)$ (Equation (3)). For instance, in the case of a single masked position (i.e., $|M| = 1$), the loss is equal to the logarithm of perplexity. More generally, it can be shown formally that by minimizing the masked language modeling objective, the model learns to approximate the conditional marginals of the language (of proteins), including the leave-one-out probabilities evaluated in perplexity (Hennigen & Kim, 2023). As a result, applying self-supervised test-time customization on x through $g \circ f$ enhances the representation of the target protein in the backbone f , leading to improved downstream performance via the fine-tuning track $h \circ f$.

B CUSTOMIZATION BEYOND MASKED LANGUAGE MODELING

In this work, we primarily focus on protein language models pretrained with masked language modeling (MLM), where a fixed proportion of randomly selected tokens (e.g., 15%) are masked for training. To date, MLM has been the dominant paradigm in protein representation learning. Nevertheless, we also provide a proof of concept showing that ProteinTTT can be applied to autoregressive and discrete

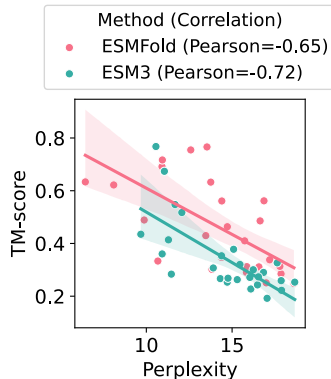


Figure A1: Quality of protein structure prediction, as measured by TM-score, correlates with perplexity of the underlying language model on the challenging targets from the CAMEO validation set. Higher TM-scores are associated with lower perplexity, indicating that better predictions are linked to lower uncertainty in the language model’s understanding of the protein sequence.

²Please note that this is an approximation of perplexity, which is computationally intractable for bidirectional models, and is often referred to as pseudo-perplexity (Lin et al., 2023; Salazar et al., 2019).

diffusion-based protein language models, with details provided in the corresponding paragraphs below. Furthermore, in Appendix I we discuss how ProteinTTT could be extended beyond protein language models.

Autoregressive customization objective. To perform single-sequence customization in an autoregressive setting (i.e., customization of ProGen2 (Nijkamp et al., 2023)), we apply a standard teacher forcing procedure (Vaswani, 2017) with a batch size of one. Specifically, each ProteinTTT step optimizes next token prediction across the whole sequence in parallel via the following loss function:

$$\mathcal{L}_{\text{AR}}(x; \theta) = \frac{1}{|x|} \sum_{i=1}^{|x|} -\log p(x_i | x_{<i}; \theta), \quad (4)$$

where x denotes a sequence of protein tokens, and $p(x_i | x_{<i}; \theta) \doteq g(f(x_{<i}; \theta))_{x_i}$ is the probability assigned by the model to the true token x_i given all preceding tokens $x_{<i}$. Here, we use the notation consistent with Equation (2).

Discrete diffusion customization objective. Recently, discrete diffusion protein language models have emerged as an extension of MLM-based protein language models. Instead of masking a fixed ratio of tokens, discrete diffusion approaches vary the masking ratio during training according to a diffusion schedule (Hsieh et al., 2025; Wang et al., 2024b;a; Campbell et al., 2024; Alamdari et al., 2023). This has been shown to improve representation learning and to enable sequence generation by starting from a fully masked sequence and gradually denoising it (Wang et al., 2024a).

In this work, we experiment with the DPLM2 Bit-based discrete diffusion model (Hsieh et al., 2025) for protein structure prediction. Interestingly, we find that using a standard MLM objective with a fixed 15% masking ratio for customization (Equation (2)) already improves performance. Exploring modifications of the customization objective tailored specifically to discrete diffusion models presents an exciting direction for future work.

C CUSTOMIZATION WITH MULTIPLE SEQUENCE ALIGNMENT (MSA)

Table A1: **ProteinTTT can be used with MSA when available.** Please see Table 2 for evaluation details.

Method	Avg. Spearman \uparrow
ESM2 (Lin et al., 2023)	0.4139
ESM2 + ProteinTTT _{MSA} (Ours)	0.4299 \pm 0.00099
MSA Transformer (Rao et al., 2021)	0.4319
MSA Transformer + ProteinTTT (Ours)	0.4326 \pm 0.00003

Customization training objective. Since many target proteins may not have homologous sequences (Rao et al., 2021) and finding such homologs may be time-consuming (Lin et al., 2023), the ProteinTTT customization objective (Equation (2)) only assumes a single target sequence for customization. However, we also extend the loss function to the case when a multiple sequence alignment (MSA) is available:

$$\mathcal{L}_{\text{MSA}}(x; \theta) = \mathbb{E}_{x' \sim p_{\text{MSA}}(x'|x)} [\mathcal{L}(x'; \theta)], \quad (5)$$

where $p_{\text{MSA}}(x'|x)$ is the distribution of sequences x' homologous to the target protein x , \mathcal{L} is the single-sequence loss function defined in Equation (2), and θ denotes the tunable parameters of the model backbone f . We refer to customization using Equation (5) as ProteinTTT_{MSA}.

Results for fitness prediction. It is known that evolutionary information is important for protein fitness prediction (Laine et al., 2019). Therefore, we demonstrate how ProteinTTT_{MSA} and ProteinTTT can enhance the performance of PLMs on the ProteinGym benchmark (Notin et al., 2024). Table A1 shows that using ProteinTTT_{MSA} with high-quality MSAs curated by Notin et al. (2024) strongly enhances the performance of ESM2, approaching that of MSA Transformer, pre-trained on MSAs. Moreover, we find that MSA Transformer slightly benefits from single-sequence customization

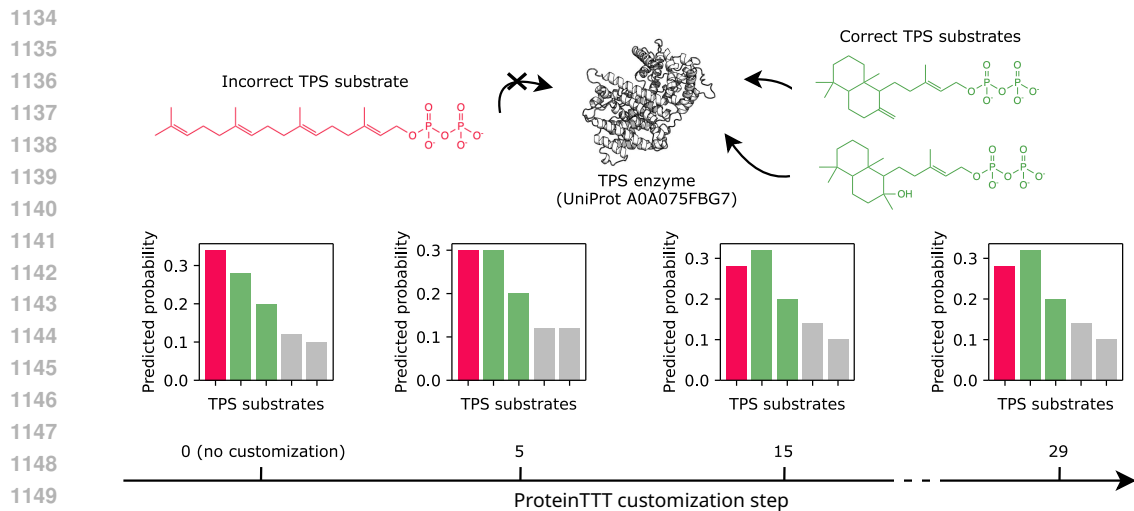


Figure A2: **Customization with ProteinTTT enables the correct substrate classification for a terpene synthase (TPS) enzyme.** With progressive customization steps of EnzymeExplorer + ProteinTTT, the probability of the initially misclassified substrate (red) decreases, while the probability of the true substrates (green) increases. The bar plots also display the predicted probabilities for other substrates with non-zero values (grey).

with ProteinTTT, while customization to whole or subsampled MSAs disrupts the performance (Table A3 in Appendix H.2). Please note that similar results were previously demonstrated in (Gordon et al., 2024) and (Alley et al., 2019) by fine-tuning protein language models on MSA, referred to as “evotuning”.

D CUSTOMIZATION FOR PROTEIN FUNCTION PREDICTION

Protein function prediction is essential for understanding biological processes and guiding bioengineering, but is challenging due to its vague definition and limited data (Yu et al., 2023; Radivojac & et al., 2013; Stärk et al., 2021; Mikhael et al., 2024; Samusevich et al., 2025). While improved structure prediction with ProteinTTT (Section 4.1) can already enhance function prediction (Song et al., 2024), we also evaluate our customization method directly on two function classification tasks: subcellular localization, predicting protein location within a cell (Stärk et al., 2021), and substrate classification for terpene synthases (TPS), enzymes producing the largest class of natural products (Christianson, 2017; Samusevich et al., 2025). Using ProteinTTT with EnzymeExplorer (Samusevich et al., 2025) for TPS detection and Light attention (Stärk et al., 2021) for subcellular localization, we achieve consistent performance gains.

Evaluation setup. For the terpene substrate classification, we use the largest available dataset of characterized TPS from Samusevich et al. (2025) and reuse the original cross-validation schema. In the case of protein localization prediction, we use a standard DeepLoc dataset (Almagro Armenteros et al., 2017) as a validation set and setHard from Stärk et al. (2021) as the test set.

Given a protein, the goal of function prediction is to correctly classify it into one of the predefined functional annotations. We assess the quality of the TPS substrate prediction using standard multi-label classification metrics used in the EnzymeExplorer paper (Samusevich et al., 2025): mean average precision (mAP) and area under the receiver operating characteristic curve (AUROC). In the case of protein localization prediction, we similarly use the classification metrics from the original paper (Stärk et al., 2021): accuracy, multi-class Matthews correlation coefficient (MCC), and F1-score.

Results. Customization with ProteinTTT improves model performance on both of the protein function prediction tasks and across all considered metrics (Table A2). Figure A2 provides a qualitative

Table A2: **Customization with ProteinTTT improves protein function prediction.** For the terpene synthase (TPS) substrate classification task, the metrics are computed on the 512 TPS sequences based on the cross-validation schema of the TPS dataset (Samusevich et al., 2025). Subcellular localization prediction performance is reported for 432 protein sequences from the setHard test set (Stärk et al., 2021). The error bars show standard deviations across five random seeds.

TPS substrate classification			Subcellular localization prediction			
Method	mAP \uparrow	AUROC \uparrow	Method	Accuracy \uparrow	MCC \uparrow	F1-score \uparrow
EnzymeExplorer (Samusevich et al., 2025)	0.805	0.948	Light attention (Stärk et al., 2021)	0.627	0.549	0.618
EnzymeExplorer + ProteinTTT (Ours)	0.811 \pm 0.0011	0.950 \pm 0.0002	Light attention + ProteinTTT (Ours)	0.634 \pm 0.004	0.557 \pm 0.005	0.627 \pm 0.004

result, where customization with ProteinTTT iteratively refines the prediction of EnzymeExplorer toward a correct TPS substrate class. We hypothesize that improvement with customization is more challenging in classification tasks, as opposed to regression problems, because a larger change in the latent space is required to shift the top-class probability.

E IMPLEMENTATION DETAILS

```

1211 1 import esm
1212 2 from proteinttt.models.esmfold import ESMFoldTTT, DEFAULT_ESMFOLD_TTT_CFG
1213 3
1214 4 # Set protein sequence
1215 5 sequence = (
1216 6     "GIHLGELGLLPSTVLAIGYFENLVNIIICESLNMLPKLEVSGKEYKKFKFTIVIPKDLLDANIKKRAKIY"
1217 7     "FKQKSLIEIEIPTSSRNYP IHIQFDENSTD DILHLYDMP TTI GGIDKAIEMFMRKGHIGKTDQQKLE"
1218 8     "ERELRNFKTTLENLIATDAFAKEMVEVIIEE"
1219 9 )
1220 10
1221 11 # Load model
1222 12 model = esm.pretrained.esmfold_v1()
1223 13 model = model.eval().cuda()
1224 14
1225 15 predict_structure(model, sequence)
1226 16 # pLDDT: 38.43025
1227 17
1228 18 # ===== ProteinTTT =====
1229 19 # Customize model to sequence
1230 20 model = ESMFoldTTT.ttt_from_pretrained(
1231 21     model, ttt_cfg=DEFAULT_ESMFOLD_TTT_CFG, esmfold_config=model.cfg
1232 22 )
1233 23 model.ttt(sequence)
1234 24 # =====
1235 25
1236 26 predict_structure(model, sequence)
1237 27 # pLDDT: 78.69619
1238 28
1239 29 # ===== ProteinTTT =====
1240 30 # Reset model to original state (after this model.ttt can be called with
1241 31 # another protein)
1242 32 model.ttt_reset()
1243 33 # =====

```

Code snippet 1: Incorporation of ProteinTTT into an ESMFold structure prediction pipeline using the proteinttt package.

```

1242 1 import torch
1243 2 import esm
1244 3 from esm.model.esm2 import ESM2
1245 4 from proteinttt.base import TTTModule
1246 5
1247 6 class ESM2TTT(TTTModule, ESM2):
1248 7     def __init__(self, ttt_cfg: TTTConfig, **kwargs):
1249 8         ESM2.__init__(self, **kwargs)
1250 9         TTTModule.__init__(self, ttt_cfg=ttt_cfg)
1251 10        self.ttt_alphabet = esm.Alphabet.from_architecture("ESM-1b")
1252 11        self.ttt_batch_converter = self.ttt_alphabet.get_batch_converter()
1253 12
1254 13    def _ttt_tokenize(self, seq: str, **kwargs):
1255 14        batch_labels, batch_strs, batch_tokens = self.ttt_batch_converter(
1256 15            [(None, seq)]
1257 16        )
1258 17        return batch_tokens
1259 18
1260 19    def _ttt_get_frozen_modules(self) -> list[torch.nn.Module]:
1261 20        return [self.embed_tokens]
1262 21
1263 22    def _ttt_mask_token(self, token: int) -> int:
1264 23        return self.ttt_alphabet.mask_idx
1265 24
1266 25    def _ttt_get_padding_token(self) -> int:
1267 26        return self.ttt_alphabet.padding_idx
1268 27
1269 28    def _ttt_token_to_str(self, token: int) -> str:
1270 29        return self.ttt_alphabet.all_toks[token]
1271 30
1272 31    def _ttt_get_all_tokens(self) -> list[int]:
1273 32        return [
1274 33            self.ttt_alphabet.tok_to_idx[t]
1275 34            for t in self.ttt_alphabet.all_toks
1276 35        ]
1277 36
1278 37    def _ttt_get_non_special_tokens(self) -> list[int]:
1279 38        return [
1280 39            self.ttt_alphabet.tok_to_idx[t]
1281 40            for t in self.ttt_alphabet.standard_toks
1282 41        ]
1283 42
1284 43    def _ttt_predict_logits(
1285 44        self, batch: torch.Tensor, start_indices: torch.Tensor = None
1286 45    ) -> torch.Tensor:
1287 46        return self(batch) ["logits"]

```

Code snippet 2: Implementation of ESM2 + ProteinTTT within the proteinttt package.

Infrastructure. All experiments with ProteinTTT are conducted on machines equipped with a single NVIDIA A100 40GB GPU, an 8-core AMD processor, and 128 GB of physical memory.

Source code. We provide a user-friendly and easily extensible PyTorch (Paszke, 2019) implementation of ProteinTTT, available as the proteinttt Python package³. We provide Code snippet 1 and Code snippet 2 in Python to demonstrate the implementation of inference and customization with ProteinTTT, respectively. Code snippet 1 demonstrates how inference with ESMFold can be enhanced with ProteinTTT by adding just a few lines of code to enable customization. Next, Code snippet 2 shows how ProteinTTT can be easily implemented for a PLM of interest by inheriting from the abstract TTTModule class. To integrate ProteinTTT within a model (e.g., ESM2), the user needs

³<https://anonymous.4open.science/r/ProteinTTT-anonymous-F585>

1296 to implement methods that define the model’s vocabulary, an interface for predicting logits, and a
1297 specification of which modules need to be fine-tuned or remain frozen. The rest, i.e., the test-time
1298 training logic itself, is implemented within the unified `TTTModule` class.
1299

1300 **Optimization.** We minimize the loss defined in Equation (2) using stochastic gradient descent
1301 (SGD) with zero momentum and zero weight decay (Ruder, 2016). While a more straightforward
1302 option might be to use the optimizer state from the final pre-training step, this approach is often
1303 impractical because the optimizer parameters are usually not provided with the pre-trained model
1304 (Hayes et al., 2024; Lin et al., 2023). Moreover, many models are pre-trained using the Adam
1305 optimizer (Kingma & Ba, 2015) or its variants (Loshchilov & Hutter, 2019). However, it was shown
1306 that Adam results in less predictable behavior of test-time training compared to the SGD optimizer,
1307 possibly due to its more exploratory behavior (Gandelsman et al., 2022).
1308

1309 **Customizing large models.** We aim for customization to be applicable on the fly, i.e., without
1310 the need for any pre-computation and on a single GPU with a minimum computational overhead.
1311 Since state-of-the-art models for many protein-oriented tasks are typically large, with up to billions
1312 of parameters, our aim presents two key challenges. First, when using pre-trained Transformers on
1313 a single GPU, even for the forward pass, the batch size is typically limited to only several samples
1314 due to the quadratic complexity of the inference (Vaswani, 2017). Second, for the backward pass,
1315 even a batch size of one is not always feasible for large models. To address the first challenge, we
1316 perform forward and backward passes through a small number of training examples and accumulate
1317 gradients to simulate updates with any batch size. We address the second challenge by employing
1318 low-rank adaptation (LoRA; Hu et al. (2021)), which in practice enables fine-tuning of any model for
1319 which a forward pass on a single sample is feasible, due to a low number of trainable parameters.
1320 Appendix H.3 details how ESMFold (Lin et al., 2023), with its 3B-parameter ESM2 backbone f , can
1321 be efficiently customized, retaining its speed advantage while enhancing performance.
1322

1323 F EXPERIMENTAL DETAILS

1324
1325 In this section, we describe the proposed benchmark suite for the three customization tasks con-
1326 sidered in this work: protein structure prediction (Appendix F.1), protein fitness prediction (Ap-
1327 pendix F.2), and protein function prediction (Appendix F.3). Each subsection describes the application
1328 of ProteinTTT to the respective models, along with details on the data, metrics, and models. Table A3
1329 additionally summarizes the hyperparameters used for the application of ProteinTTT to individual
1330 models.
1331

1332 F.1 PROTEIN STRUCTURE PREDICTION

1333 F.1.1 DATASETS

1334
1335 **CAMEO dataset.** To evaluate the capabilities of ProteinTTT on protein structure prediction, we
1336 employ the CAMEO validation and test sets as described in Lin et al. (2023). Specifically, the
1337 validation set was obtained by querying the CAMEO (Continuous Automated Model Evaluation) web
1338 server⁴ (Robin et al., 2021) for entries between August 2021 and January 2022, while the CAMEO
1339 test set consists of entries from April 1, 2022, to June 25, 2022. Most of the entries in the CAMEO
1340 sets are predicted with high accuracy and confidence (Lin et al., 2023). Therefore, we subselect the
1341 challenging validation and test sets where customization with ProteinTTT is suitable.

1342 Specifically, we apply two standard criteria: (1) preserving entries with ESMFold pLDDT scores
1343 below 70 to filter out high-confidence predictions (Jumper et al., 2021), and (2) selecting entries
1344 with ESM2 perplexity scores greater than or equal to 6, ensuring that the predictions are challenging
1345 due to poor sequence understanding rather than other factors. Additionally, most structures with
1346 perplexity scores below 6 are already associated with high-confidence predictions (Figure S5 in Lin
1347 et al. (2023)). After filtering, the resulting challenging validation and test sets consist of 27 (out of
1348 378) and 18 (out of 194) targets, respectively.
1349

⁴<https://www.cameo3d.org/modeling>

1350 F.1.2 METRICS

1351
1352 To assess the quality of the predicted protein structures with respect to the ground truth structures, we
1353 use two standard metrics averaged across the test dataset: TM-score (Zhang & Skolnick, 2004) and
1354 LDDT (Mariani et al., 2013).

1355
1356 **TM-score.** The TM-score (Template Modeling score) is a metric used to assess the quality of the
1357 global 3D alignment between the predicted and target protein structures. It evaluates the structural
1358 similarity by comparing the distance between corresponding residues after superposition. The
1359 TM-score ranges from 0 to 1, where higher values indicate better alignment.

1360
1361 **LDDT.** The Local Distance Difference Test (LDDT) is an alignment-free metric used to assess the
1362 accuracy of predicted protein structures. Unlike global metrics, LDDT focuses on local structural
1363 differences by measuring the deviation in distances between atom pairs in the predicted structure
1364 compared to the target structure. It is particularly useful for evaluating the accuracy of local regions,
1365 such as secondary structure elements. LDDT scores range from 0 to 100, with higher values indicating
1366 better local structural agreement.

1367 F.1.3 MODELS

1368
1369 **ESMFold.** The ESMFold architecture comprises two key components: a protein language model,
1370 ESM2, which, given a protein sequence, generates embeddings for individual amino acids, and a
1371 folding block that, using these embeddings and the sequence, predicts the protein 3D structure along
1372 with per-amino-acid confidence scores, known as pLDDT scores. In our experiments, we use the
1373 `esmfold_v0` model from the publicly available ESMFold checkpoints⁵. Please note that we use
1374 `esmfold_v0` and not `esmfold_v1` to avoid data leakage with respect to the CAMEO test set.

1375
1376 **ESMFold + ProteinTTT.** Since the ESM2 backbone of ESMFold was pre-trained in a self-
1377 supervised masked modeling regime, the application of ProteinTTT to ESMFold is straightforward.
1378 We treat ESM2 as the backbone f , the language modeling head predicting amino acid classes from
1379 their embeddings as the self-supervised head g , and the folding trunk along with the structure modules
1380 as the downstream task head h . After each ProteinTTT step, we run $h \circ f$ to compute the pLDDT
1381 scores, which allows us to estimate the optimal number of customization steps for each protein based
1382 on the highest pLDDT score.

1383 Since the backbone f is given by the ESM2 model containing 3 billion parameters, we apply LoRA
1384 (Hu et al., 2021) to all matrices involved in self-attention. This enables fine-tuning ESMFold +
1385 ProteinTTT on a single GPU.

1386
1387 **ESMFold + ME.** Since ESMFold is a regression model, it only predicts one solution and does not
1388 have a straightforward mechanism for sampling multiple structure predictions. Nevertheless, the
1389 authors of ESMFold propose a way to sample multiple candidates (Section A.3.2 in Lin et al. (2023)).
1390 To sample more predictions, the masking prediction (ME) method randomly masks 15% (same ratio
1391 as during masked language modeling pre-training) of the amino acids before passing them to the
1392 language model. Selecting the solution with the highest pLDDT may lead to improved predicted
1393 structure. Since sampling multiple solutions with ESMFold + ME and selecting the best one via
1394 pLDDT is analogous to ESMFold + ProteinTTT, we employ the former as a baseline, running the
1395 method for the same number of steps.

1396
1397 **HelixFold-Single.** HelixFold-Single is an MSA-free protein structure prediction model that com-
1398 bines representations from a pretrained protein language model with adapted AlphaFold2 geometric
1399 modules (EvoformerS and Structure) to directly predict atomic coordinates (Fang et al., 2023). We
1400 use the official implementation⁶

1401 ⁵[https://github.com/facebookresearch/esm/blob/main/esm/esmfold/v1/](https://github.com/facebookresearch/esm/blob/main/esm/esmfold/v1/pretrained.py)
1402 [pretrained.py](https://github.com/facebookresearch/esm/blob/main/esm/esmfold/v1/pretrained.py)

1403 ⁶[https://github.com/PaddlePaddle/PaddleHelix/tree/dev/apps/protein_](https://github.com/PaddlePaddle/PaddleHelix/tree/dev/apps/protein_folding/helixfold-single)
[folding/helixfold-single](https://github.com/PaddlePaddle/PaddleHelix/tree/dev/apps/protein_folding/helixfold-single)

1404 **HelixFold-Single + ProteinTTT.** HelixFold-Single shares the main concept with ESMFold, and
 1405 we combine it with ProteinTTT in the same way as in ESMFold + ProteinTTT.

1407 **DPLM2 Bit-based.** The DPLM2 Bit-based discrete diffusion protein language model (Hsieh
 1408 et al., 2025) extends DPLM2 by using bit-wise discrete modeling to enhance structure generation
 1409 capabilities (Wang et al., 2024b). DPLM2 is a multi-modal model that jointly models protein
 1410 sequences and discretized structural tokens within a single discrete diffusion framework. In this work,
 1411 we evaluate DPLM2 Bit-based on the task of structure prediction. Structure prediction is performed
 1412 by initializing the structural tokens with masks and gradually denoising them based on the sequential
 1413 tokens. We use the official implementation⁷ with the standard 650M-parameter model, 100 denoising
 1414 steps, and the denoising strategy set to `annealing@1.1:0.1`.

1415 **DPLM2 Bit-based + ProteinTTT.** To apply ProteinTTT to DPLM2 Bit-based, we use the standard
 1416 masked language modeling objective (Equation (2)). See Appendix B for further discussion. Please
 1417 also note that we do not use confidence function c with DPLM2 Bit-based as it does not implement
 1418 pLDDT or any other confidence function for protein structure prediction.

1420 **ESM3.** Unlike ESMFold, ESM3 is a fully multiple-track, BERT-like model (Devlin, 2018), pre-
 1421 trained to unmask both protein sequence and structure tokens simultaneously (along with the function
 1422 tokens). The structure tokens in ESM3 are generated via a separately pre-trained VQ-VAE (Razavi
 1423 et al., 2019) operating on the protein geometry. In our experiments, we use the smallest, publicly
 1424 available version of the ESM3 model (`ESM3_sm_open_v0`)⁸.

1425 **ESM3 + ProteinTTT.** We treat the Transformer encoder of ESM3 as f , the language modeling
 1426 head decoding amino acid classes as g , and the VQ-VAE decoder, which maps structure tokens to the
 1427 3D protein structure, as h . During the customization steps, we train the model to unmask a protein
 1428 sequence while keeping the structural track fully padded. During the inference, we provide the model
 1429 with a protein sequence and run it to unmask the structural tokens, which are subsequently decoded
 1430 with the VQ-VAE decoder. After each customization step, we run $h \circ f$ to compute the pLDDT
 1431 scores, which allows us to estimate the optimal number of customization steps for each protein based
 1432 on the highest pLDDT score. We choose the optimal hyperparameters by maximizing the difference
 1433 in TM-score after and before applying ProteinTTT across the validation dataset.

1434 Despite the fact that the model contains 1.4 billion parameters, even without using LoRA, ESM3 +
 1435 ProteinTTT can be fine-tuned on a single NVIDIA A100 GPU. Therefore, we do not employ LoRA
 1436 for fine-tuning ESM3, while this can also be possible.

1437 **ESM3 + CoT.** To improve the generalization and protein-specific performance of ESM3, the
 1438 original ESM3 paper employs a chain of thought (CoT) procedure. The procedure unfolds in n steps
 1439 as follows. At each step, $1/n$ of the masked tokens with the lowest entropy after softmax on logits
 1440 are unmasked. Then, the partially unmasked sequence is fed back into the model, and the process
 1441 repeats until the entire sequence is unmasked. In our experiments, we set $n = 8$, which is the default
 1442 value provided in the official GitHub repository.

1443 F.2 PROTEIN FITNESS PREDICTION

1444 F.2.1 DATASETS

1445 **ProteinGym.** ProteinGym⁹ is the standard benchmark for protein fitness prediction (Notin et al.,
 1446 2024). The latest, second version of the dataset includes 217 deep mutation scanning experiments
 1447 (DMSs) across different proteins. We focus on the well-established zero-shot setup of the benchmark
 1448 and do not experiment with the supervised setup, as it has not yet been fully incorporated into the
 1449 official codebase at the time of this study. In total, the dataset contains 2.5M mutants with annotated

1450 ⁷<https://github.com/bytedance/dplm>

1451 ⁸<https://github.com/evolutionaryscale/esm>

1452 ⁹<https://github.com/OATML-Markslab/ProteinGym>

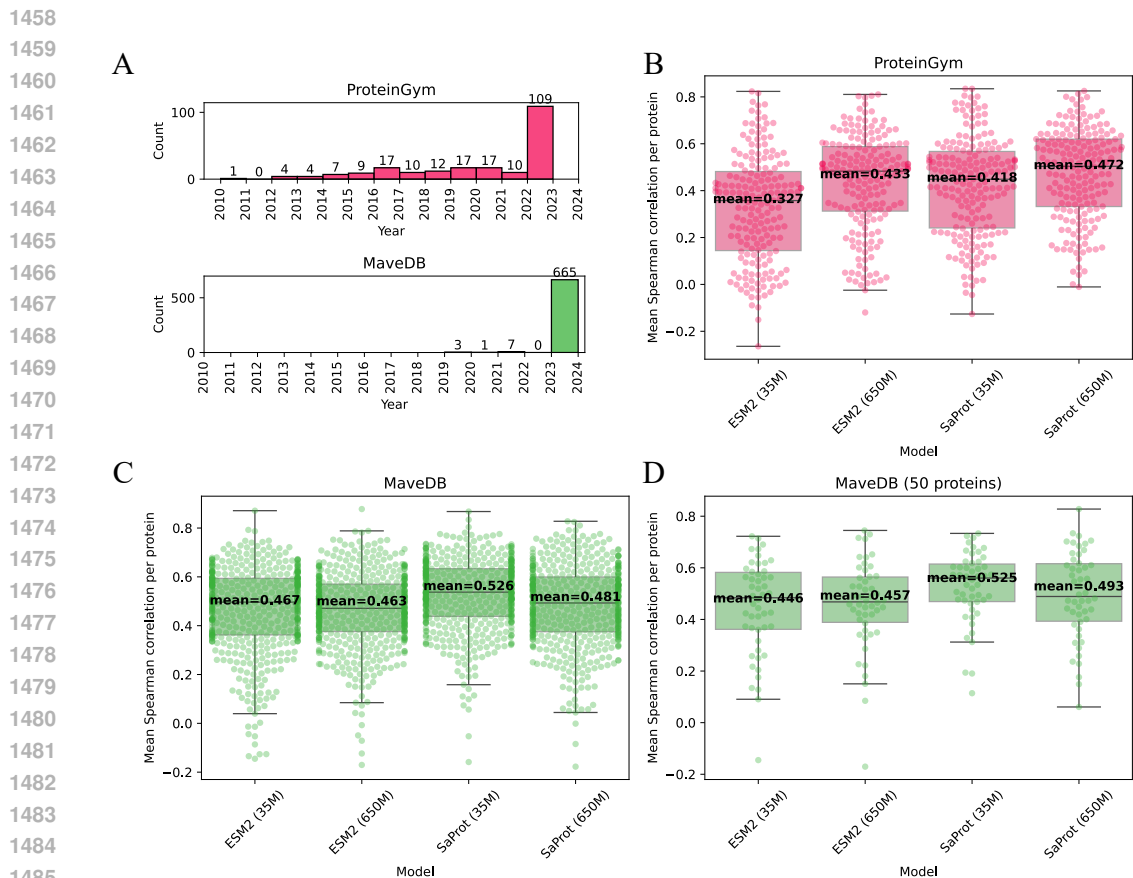


Figure A3: **Comparison of the standard ProteinGym dataset with the MaveDB dataset constructed in this work.** A) MaveDB, mined from Esposito et al. (2019), includes novel assays even after filtering to ensure distinct proteins from the comprehensive ProteinGym dataset. This is largely because most MaveDB assays post-filtering date to 2024, whereas the latest assays in ProteinGym date to 2023. B, C, D) MaveDB is of sufficient quality for model evaluation. Representative baselines, ESM2 and SaProt with both 35 million and 650 million parameters, evaluated on ProteinGym generalize effectively to MaveDB, following a similar distribution of predictions. Panel D illustrates the random subset of 50 proteins used for hyperparameter tuning for fitness prediction. Each point in the plots represents one protein and shows the Spearman correlation averaged across all assays corresponding to the protein (typically one assay per protein). The box plots standardly depict quartiles, medians, and outliers.

ground-truth fitness. Since ProteinGym does not contain a data split for the zero-shot setup, employed in this work, we use the whole dataset as the test set.

MaveDB dataset. To establish a validation set disjoint from ProteinGym (Notin et al., 2024), we mined MaveDB¹⁰ (Esposito et al., 2019). As of August 1, 2024, the database contains 1178 Multiplexed Assays of Variant Effects (MAVEs), where each assay corresponds to a single protein, measuring the experimental fitness of its variants. We applied quality control filters to remove potentially noisy data. Specifically, we ensured that the UniProt identifier (Consortium, 2023) is valid and has a predicted structure available in the AlphaFold DB (Varadi et al., 2022). We also excluded assays with fewer than 100 variants, as well as those where at least one mutation had a wrongly annotated wild type or where most mutations failed during parsing. Additionally, to ensure no overlap between datasets, we removed any assays whose UniProt identifier matched with those in ProteinGym, ensuring that the validation and test sets contain different proteins.

¹⁰<https://www.mavedb.org>

The described methodology resulted in the MaveDB dataset comprising 676 assays (out of 1178 in the entire MaveDB) with experimental fitness annotations. This corresponds to 483 unique protein sequences and 867 thousand mutations in total. The large size of the dataset, despite the comprehensiveness of ProteinGym containing 217 assays, can be attributed to the fact that many assays in MaveDB were released after the ProteinGym construction (Figure A3A). To ensure the quality of the constructed MaveDB dataset, we validated that representative baselines from ProteinGym generalize to the new assays, following similar distributions of predictions (Figure A3B,C). Finally, for efficiently tuning hyperparameters for fitness prediction models, we sampled 50 proteins (Figure A3D), corresponding to 83 assays comprising 134 thousand variants.

F.2.2 METRICS

Protein fitness labels are not standardized and can vary across different proteins. Nevertheless, the ranking of mutations for a single protein, as defined by fitness labels, can be used to assess the mutation scoring capabilities of machine learning models. As a result, Spearman correlation is a standard metric for evaluation.

Spearman by phenotype. When computing Spearman correlations, we follow the evaluation protocol proposed in ProteinGym (Notin et al., 2024). First, for each protein, we compute Spearman correlation scores between the predicted ranks of mutations and their corresponding labels. Then, we average the scores across five categories of assayed phenotypes, measuring the effects of mutations: catalytic activity (“Activity”), binding affinity to a target (“Binding”), protein expression levels in a cell (“Expression”), organism growth rate (“Organismal Fitness”), and protein thermostability (“Stability”).

Avg. Spearman. We refer to the mean score across the five phenotype categories as “Avg. Spearman”. We report the “Avg. Spearman” metric as the mean and standard deviation across five random seeds (Table 2, Table A4).

Spearman by MSA Depth. Following (Notin et al., 2024), we split the performance by the depth of available multiple sequence alignment (MSA), i.e., the number of homologous sequences available, as provided in ProteinGym: “Low depth”, “Medium depth”, and “High depth”, and report the Spearman correlation for each subset individually (Table A4). Specifically, the MSA depth categories in ProteinGym are determined using the following thresholds from Hopf et al. (2017): “Low” is defined as $N_{eff}/L < 1$, “Medium” as $1 < N_{eff}/L < 100$, and “High” as $N_{eff}/L > 100$, where N_{eff} represents the normalized number of effective sequences in the MSA, and L is the sequence length covered in the MSA.

F.2.3 MODELS

ESM2. The ESM2 model is a bidirectional, BERT-like (Devlin, 2018) Transformer trained on millions of protein sequences using masked modeling (Lin et al., 2023). The goal of protein fitness prediction is to predict the effects of mutations, and PLMs are often adapted to this task using zero-shot transfer via log odds ratio (Notin et al., 2024; Meier et al., 2021). Specifically, for a given single- or multi-point mutation, where certain amino acids T are substituted from x_i to x_i^m for each $i \in T$, the fitness prediction via the log odds ratio is defined as:

$$\sum_{i \in T} \left(\log p(x_i^m | x_{\setminus i}) - \log p(x_i | x_{\setminus i}) \right), \quad (6)$$

where the sum iterates over mutated positions $i \in T$ with $p(x_i^m | x_{\setminus i})$ and $p(x_i | x_{\setminus i})$ denoting the predicted probabilities of the mutated amino acid and the original one (i.e., wild type), respectively. The conditionals $x_{\setminus i}$ indicate that the input sequence to the model has the position i masked. In this setup, the native (unmutated) sequence, where $T = \emptyset$, has a predicted fitness of 0. Mutations with negative values represent favorable mutations, while positive values correspond to disruptive mutations. We follow the ProteinGym benchmark and use this formula (Notin et al., 2024) to evaluate the fitness prediction capabilities of ESM2. We use the implementation of ESM2 from ProteinGym.

ESM2 + ProteinTTT. ESM2 can be straightforwardly customized with ProteinTTT. Specifically, we treat the Transformer encoder as the backbone f , and the language modeling head, which projects token embeddings to amino acid probabilities, as the pre-training head g . The log odds ratio given by Equation (6) serves as the task-specific head h , which in this case involves the pre-training head g that predicts log probabilities. Overall, we apply ProteinTTT to the pre-trained ESM2 model and, after a pre-defined number of self-supervised fine-tuning steps, score mutations using Equation (6). During customization, we fine-tune all parameters in $g \circ f$ end-to-end except for token and position embeddings. When evaluating ESM2 + ProteinTTT_{MSA}, we use the MSAs curated by the authors of ProteinGym (Notin et al., 2024).

SaProt. We also experiment with a structure-aware protein language model, SaProt (Su et al., 2023). SaProt builds off the ESM2 model but incorporates structural information from predicted protein structures. Specifically, SaProt uses the same Transformer architecture but expands its vocabulary by combining the 20 standard amino acid tokens with 20 structural tokens from the 3Di vocabulary, increasing the total alphabet size to 400. The 3Di tokens capture the geometry of the protein backbone and are generated using VQ-VAE (Razavi et al., 2019), which projects continuous geometric information into discrete tokens and was trained as part of the Foldseek method (van Kempen et al., 2022).

Since SaProt is also a protein language model, it also uses Equation (6) to score variants. However, please note that SaProt, as implemented in ProteinGym (Notin et al., 2024), uses a slightly different version of the log odds ratio. In SaProt, the conditions in the log probabilities in Equation (6) are replaced with $x_{\setminus T}$ instead of $x_{\setminus i}$, not assuming the independence of substitutions. During customization with ProteinTTT, we only mask sequential information and leave the structural part of the tokens unchanged, reflecting the original pre-training setup. We use the implementation of SaProt from ProteinGym⁹.

SaProt + ProteinTTT. Since the architecture of SaProt is based on ESM2, the ProteinTTT components f , g , and h remain the same. It means that customization can be applied to the model in the same way as in the case of ESM2 + ProteinTTT discussed above.

ProSST. We experiment with the state-of-the-art fitness predictor, ProSST (Li et al., 2024). ProSST primarily improves upon SaProt (Su et al., 2023) by incorporating a larger vocabulary of structural tokens and employing disentangled attention mechanisms. Instead of relying on the 3Di alphabet optimized for protein structure search with Foldseek (van Kempen et al., 2022), Li et al. (2024) pre-train a new autoencoder to denoise corrupted protein backbones and cluster the resulting latent space using the K -means algorithm (Lloyd, 1982). Notably, optimal performance for fitness prediction is achieved with $K = 2048$ tokens, compared to just 20 in the 3Di vocabulary used by SaProt. We adopt this model in our experiments. Additionally, disentangled attention in ProSST enhances information propagation between sequence and structure within its Transformer blocks, further improving prediction performance. The model has 110M parameters in total.

ProSST, similarly to ESM2 and SaProt, is pre-trained using masked language modeling applied to protein sequence tokens. To score mutations on the ProteinGym benchmark (Notin et al., 2024), ProSST also uses the log-odds ratio, but in a slightly different way compared to ESM2 and SaProt. Specifically, ProSST performs a single forward pass to predict log probabilities, which are then used to score all mutations. Formally, this approach modifies the log probability condition in Equation (6), replacing $x_{\setminus i}$ with x .

ProSST + ProteinTTT. Similarly to ESM2 and SaProt, we treat the Transformer encoder in ProSST as the backbone f , the masked language modeling head as the pre-training head g , and the log-odds ratio formula as the task-specific head h .

ProGen2. For fitness prediction, we additionally experiment with one of the major autoregressive protein language models, ProGen2 (Nijkamp et al., 2023). Specifically, we experiment with ProGen of two sizes: ProGen2-small (151M parameters) and ProGen2-large (2.7B parameters). We obtain the pre-trained weights from the official GitHub repository¹¹. For ProGen2-large inference, we use floating-point 16 precision for computational efficiency.

¹¹<https://github.com/salesforce/progen/tree/main/progen2>

1620 **ProGen2 + ProteinTTT.** To demonstrate the applicability of ProteinTTT in an autoregressive
1621 setting, we apply it to the ProGen2 (Nijkamp et al., 2023) language model. To perform the cus-
1622 tomization, we use the standard next-token prediction objective on a single target protein, following
1623 Equation (4). Please see Appendix B for details.

1624
1625 **MSA Transformer.** Finally, we experiment with MSA Transformer for fitness prediction (Rao et al.,
1626 2021). Similar to ESM2 (Lin et al., 2023), MSA Transformer is pre-trained on large protein sequence
1627 datasets; however, it is trained on multiple sequence alignments (MSAs) rather than individual
1628 sequences.

1629 Since MSA Transformer is also a protein language model, it can be used for fitness prediction
1630 in the same way as ESM2, as discussed above, by computing the log-odds ratio over the first
1631 sequence in the MSA in this case. We reproduce the results of MSA Transformer on the ProteinGym
1632 benchmark with two modifications: (1) we sample a weighted subset of 32 sequences from each MSA
1633 instead of 400, and (2) we use only one random seed instead of five for ensembling. These changes
1634 significantly reduce computational time while also slightly improving performance compared to the
1635 results reported in ProteinGym. This improvement may be explained by the fact that the performance
1636 of MSA Transformer saturates with increasing MSA input size (Figure 4 in Rao et al. (2021)).

1637 **MSA Transformer + ProteinTTT.** We experiment with customizing MSA Transformer to MSA
1638 subsamples of varying sizes, ranging from a single target sequence (i.e., customization via Equation (2)
1639 with ProteinTTT) to the full MSA subset of 32 sequences (i.e., customization via Equation (5) with
1640 ProteinTTT_{MSA}). We observe that applying ProteinTTT_{MSA} to MSA Transformer with a batch size
1641 of 32 disrupts performance, while reducing the input MSA subsample size mitigates this effect.
1642 Ultimately, MSA Transformer + ProteinTTT results in a slight performance improvement.

1644 F.3 PROTEIN FUNCTION PREDICTION

1645 F.3.1 DATASETS

1646 **TPS dataset.** For the evaluation of terpene substrate classification, we use the largest available
1647 dataset of characterized TPS enzymes from Samusevich et al. (2025) and repurpose the original
1648 5-fold cross-validation schema. We focus on the most challenging TPS sequences, defined as those
1649 predicted by the TPS detector, proposed by the dataset authors, with confidence scores below 0.8.
1650 This filtering results in 104, 98, 113, 100, 97 examples in the individual folds.

1651 **setHard.** For the test evaluation of subcellular location prediction, we use the setHard dataset
1652 constructed by Stärk et al. (2021). The dataset was redundancy-reduced, both within itself and
1653 relative to all proteins in DeepLoc (Almagro Armenteros et al. (2017); next paragraph), a standard
1654 dataset used for training and validating machine learning models. The setHard dataset contains 490
1655 protein sequences, each annotated with one of ten subcellular location classes, such as “Cytoplasm”
1656 or “Nucleus”. Since we use ESM-1b (Rives et al., 2021) in our experiments with the dataset, we
1657 further filter the data to 432 sequences that do not exceed a length of 1022 amino acids. This step,
1658 consistent with Stärk et al. (2021), ensures that ESM-1b can generate embeddings for all proteins.

1659 **DeepLoc.** For hyperparameter tuning in the subcellular location prediction task, we use the test
1660 set from the DeepLoc dataset (Almagro Armenteros et al., 2017). Similar to setHard, DeepLoc
1661 assigns labels from one of ten subcellular location classes. The dataset contains 2768 proteins,
1662 which we further filter to 2457 sequences that do not exceed a length of 1022 amino acids, ensuring
1663 compatibility with the embedding capabilities of ESM-1b. Since setHard was constructed to be
1664 independent of DeepLoc, setHard provides a leakage-free source of data for validation.

1665 F.3.2 METRICS

1666 **mAP, AUROC.** The TPS substrate prediction problem is a 12-class multi-label classification task
1667 over possible TPS substrates. Therefore, we assess the quality of the predictions using standard
1668 multi-label classification metrics such as mean average precision (mAP) and area under the receiver
1669 operating characteristic curve (AUROC) averaged across individual classes. These metrics were
1670 used in the original EnzymeExplorer paper (Samusevich et al., 2025). We report the performance by
1671
1672
1673

1674 averaging the metric values concatenated across all validation folds from the 5-fold cross-validation
1675 schema.
1676

1677 **Accuracy, MCC, F1-score.** To evaluate the performance of subcellular location prediction methods,
1678 we use standard classification metrics as employed in Stärk et al. (2021). Accuracy standardly
1679 measures the ratio of correctly classified proteins, while Matthew’s correlation coefficient for multiple
1680 classes (MCC) serves as an alternative to the Pearson correlation coefficient for classification tasks
1681 (Gorodkin, 2004). The F1-score, the harmonic mean of precision and recall, evaluates performance
1682 from a retrieval perspective, balancing the trade-off between false positives and false negatives.
1683

1684 F.3.3 MODELS

1685 **EnzymeExplorer.** EnzymeExplorer is a state-of-the-art method for the classification of terpene
1686 synthase (TPS) substrates (Samusevich et al., 2025). The model consists of two parallel tracks. Given
1687 a protein sequence, EnzymeExplorer first computes its ESM-1v embedding (Meier et al., 2021) and
1688 a vector of similarities to the functional domains of proteins from the training dataset, based on
1689 unsupervised domain segmentation of AlphaFold2-predicted structures (Jumper et al., 2021). The
1690 ESM-1v embedding and the similarity vector are then concatenated and processed by a separately
1691 trained random forest, which predicts TPS substrate class probabilities.

1692 In our experiments, we use the “PLM only” version of the model, which leverages only ESM-1v
1693 embeddings. This version exhibits a minor performance decrease compared to the full model but
1694 exactly follows a Y-shaped architecture, allowing us to validate the effectiveness of ProteinTTT for
1695 predicting TPS substrates. We use the implementation of EnzymeExplorer available at the official
1696 GitHub page¹².
1697

1698 **EnzymeExplorer + ProteinTTT.** When applying ProteinTTT to EnzymeExplorer, we treat the
1699 frozen ESM-1v model as a backbone f , its language modeling head as a self-supervised head g , and
1700 the random forest classifying TPS substrates as a downstream supervised head h .
1701

1702 **Light Attention.** We use Light attention (Stärk et al., 2021) as a representative baseline for
1703 subcellular location prediction. Light attention leverages protein embeddings from a language model,
1704 which in our case is ESM-1b (Rives et al., 2021). The model processes per-residue embeddings via a
1705 softmax-weighted aggregation mechanism, referred to as light attention, which operates with linear
1706 complexity relative to sequence length and enables richer aggregation of per-residue information, as
1707 opposed to standard mean pooling. We re-train the model using ESM-1b embeddings on the DeepLoc
1708 dataset (Almagro Armenteros et al., 2017) using the code from the official GitHub page¹³.

1709 **Light attention + ProteinTTT.** When applying ProteinTTT to Light attention, we treat the frozen
1710 ESM-1b as the backbone f , the language modeling head of ESM-1b as the self-supervised head g ,
1711 and the Light attention block as the fine-tuning head h .
1712

1713 G CASE STUDY DETAILS

1714 G.1 MODELING ANTIBODY-ANTIGEN LOOPS

1715
1716 We download the SAbDab dataset from the official website¹⁴(Dunbar et al., 2014). We apply
1717 ProteinTTT to targets with low-confidence ESMFold predictions (pLDDT < 70) and remove se-
1718 quences longer than 400 residues due to GPU memory limitations. This results in a final set of
1719 175 antibody and 814 antigen chains. We predict the full structures using ESMFold+ProteinTTT
1720 (with the same hyperparameters tuned on the CAMEO validation set specified in Table A3) and
1721 compute LDDT scores against the corresponding PDB structures to assess local errors, which are
1722 particularly relevant for loop regions. For antibodies, we evaluate the complete structures, while
1723 for complementarity-determining regions (CDRs), we extract the CDR substructures as annotated
1724 in SAbDab according to Chothia numbering (Chothia & Lesk, 1987) and calculate LDDT on these
1725 regions.

1726 ¹²<https://github.com/pluskal-lab/EnzymeExplorer>

1727 ¹³<https://github.com/HannesStark/protein-localization>

¹⁴<https://opig.stats.ox.ac.uk/webapps/sabdab-sabpred/sabdab>

Table A3: **Hyperparameters used for adapting ProteinTTT to individual models.** The optimal hyperparameters were estimated using validation datasets corresponding to each of the considered tasks: *Fitness prediction*, *Structure prediction*, and *Function prediction*. Comma-separated lists show the values used for hyperparameter grid search, while the final values selected for computing the test results are highlighted in **bold**. Low-rank adaptation (LoRA) was only used with ESMFold, containing 3 billion parameters in the ESM2 backbone. Please note that we did not tune the number of customization steps, as adjusting the learning rate and batch size effectively controls the expected performance under the fixed number of steps, as shown in Figure A10. Therefore, we used 30 steps in all our experiments. The only exception was ESM3 + ProteinTTT, where the number of steps was set to 50 during initial experiments with different models/tasks conducted in parallel before standardizing the number of steps to 30. Bidirectional methods marked with an asterisk (“*”) used a slightly different calculation of the loss function. Specifically, the loss was propagated over all tokens, including special and non-masking tokens, while averaging the loss across all tokens simultaneously, rather than first averaging over sequences. This approach was used in the early stages of development, and we provide it in our codebase via `loss_kind = "unnormalized_cross_entropy"`. Please note that MSA Transformer always uses 1 MSA in a batch and the “Batch size” represents the number of sequences in this MSA with the target sequence always present as the first one.

	Learning rate	Batch size	Grad. acc. steps	Steps (Conf. func. c)	LoRA rank r	LoRA α
<i>Fitness prediction</i>						
ESM2 (35M) + ProteinTTT *	4e-5, 4e-4 , 4e-3	4	4, 8, 16 , 32, 64	30	-	-
ESM2 (650M) + ProteinTTT *	4e-5 , 4e-4, 4e-3	4	4, 8, 16 , 32	30	-	-
ProGen2-small (151M) + ProteinTTT	4e-5 , 4e-4 , 4e-3	4	4	1, 5, 10, 15 , 20	-	-
SaProt (35M) + ProteinTTT *	4e-5, 4e-4 , 4e-3	4	4, 8, 16, 32	30	-	-
SaProt (650M) + ProteinTTT *	4e-5 , 4e-4, 4e-3	2, 4	4, 8, 16 , 32	30	-	-
ProSST (K=2048) + ProteinTTT *	1e-5 , 4e-5, 4e-4, 4e-3	4	4, 8, 16, 32	30	-	-
ESM2 (650M) + ProteinTTT _{MSA} *	4e-6, 1e-5, 4e-5, 4e-4, 4e-3	4	2, 4	50, 100	-	-
MSA Transformer + ProteinTTT	1e-6, 3e-6, 1e-5, 3e-5, 1e-4	1, 4, 8, 16, 32	1, 2, 4, 8	30	-	-
<i>Structure prediction</i>						
ESM3 + ProteinTTT	1e-4, 4e-4, 1e-3	2	1, 4, 16	50 (pLDDT)	-	-
DPLM2 Bit-based + ProteinTTT	4e-6 , 4e-5 , 4e-4 , 4e-3	2, 4, 8	2, 4, 8	10	-	-
HelixFold-Single + ProteinTTT	4e-4 , 1e-3	4, 8, 16	1	30 (pLDDT)	-	-
ESMFold + ProteinTTT	4e-4	4	4, 8, 32, 64	30 (pLDDT)	4, 8, 32	8, 16, 32
<i>Function prediction</i>						
EnzymeExplorer + ProteinTTT	4e-4 , 1e-3	2	2, 4, 8	30	-	-
Light attention + ProteinTTT	4e-4, 1e-3, 3e-3	2	2, 4	30	-	-

G.2 EXPANDING KNOWN STRUCTURES OF VIRAL PROTEINS

We use BFVD version `archived/2023_02_v2`¹⁵. This version contains maximum-pLDDT structures from predictions generated by two strategies: (i) ColabFold (Mirdita et al., 2022) with MSAs constructed using Logan (Chikhi et al., 2024), and (ii) ColabFold with 12 additional recycle steps and MSAs constructed using Logan. In Figure 5, we also report pLDDT values for BFVD version `archived/2023_02_v1`, where structures are simply obtained from ColabFold with MSAs from Logan, i.e., strategy (i). We re-predict structures using ESMFold and ESMFold+ProteinTTT for sequences with length < 450 due to GPU memory constraints. We use the same hyperparameters tuned on the CAMEO validation set, as specified in Table A3, with the exception of 20 instead of 30 steps for computational efficiency.

H EXTENDED RESULTS

In this section, we provide additional results on test sets (Appendix H.1), discuss validation performance (Appendix H.2), and analyze the runtime performance of customization (Appendix H.3).

H.1 DETAILED TEST PERFORMANCE

In this section, we provide details on the test performance. Specifically, Table A4 shows that customization with ProteinTTT primarily enhances performance on challenging targets, characterized by a low number of similar proteins in sequence databases, as measured by MSA depth. Additionally, we provide a qualitative example illustrating how ProteinTTT substantially improves the correlation

¹⁵<https://bfvd.steineggerlab.workers.dev>

1782 between ESM2-predicted fitness and ground-truth stability by better identifying disruptive mutations
1783 in the protein core (Figure A5).

1784 Next, Figure A6 shows the distribution of ProteinTTT effects: in many cases, customization has
1785 minimal impact on performance; often, it leads to substantial improvements; and in rare cases,
1786 customization results in a decrease in performance. This positions ProteinTTT as a method for
1787 enhancing prediction accuracy, while a comprehensive analysis of its failure modes remains an
1788 important direction for future research. While we demonstrate these effects using a protein folding
1789 example, we observe a similar distribution of ProteinTTT impact across the tasks.

1790 We also observe that the overall trend of customization with ProteinTTT generally leads to improved
1791 performance, with robust consistency across random seeds. However, the progression of the per-
1792 formance curve can be rugged, particularly in classification tasks, where substantial changes in the
1793 underlying representations are required to shift the top-predicted class in the discrete probability
1794 distribution (Figure A12).

1796 H.2 VALIDATION PERFORMANCE

1798 This section discusses the performance of ProteinTTT on validation data. Table A5 illustrates the
1799 validation performance of the tested methods for fitness prediction on our newly constructed MaveDB
1800 dataset. ProteinTTT enhances the performance of all the methods.

1801 Next, we discuss the hyperparameter optimization. Table A3 provides the grid of hyperparameters
1802 explored for each model and its size, as well as specifies the optimal hyperparameters suitable
1803 for downstream applications. Figure A10 demonstrates the trend of hyperparameter tuning with
1804 optimal hyperparameter combination balancing underfitting and overfitting to a single target protein.
1805 While most of reasonable hyperparameter configurations lead to overall improvements when using
1806 customization with ProteinTTT, poorly chosen hyperparameters can have detrimental effects due
1807 to rapid overfitting. However, with a reliable predicted confidence measure, such as pLDDT, the
1808 appropriate customization step for each protein can be selected to mitigate overfitting. Figure A11
1809 demonstrates that when using ESM3 + ProteinTTT with pLDDT-based step selection for protein
1810 structure prediction, all hyperparameter configurations result in improved performance compared to
1811 the base ESM3 model.

1812 H.3 RUNTIME PERFORMANCE

1814 In this section, we demonstrate that customization with ProteinTTT can be done efficiently, with an
1815 acceptable computational overhead. Specifically, we show that ESMFold, known for being a faster
1816 alternative to more performant methods such as AlphaFold2 (Jumper et al., 2021) or AlphaFold3
1817 (Abramson et al., 2024), still remains in the category of lightweight methods even with ProteinTTT
1818 customization (Figure A4).

1819 This observation highlights the practical utility of ProteinTTT. For example, ESMFold enabled
1820 structural characterization of large metagenomics data (>617 million metagenomic sequences), which
1821 would be infeasible with AlphaFold2 (Lin et al., 2023). Nevertheless, the original ESMFold has
1822 high confidence predictions only for 36% of sequences from the metagenomic database, while
1823 the other 392 million sequences remain with low or medium confidence predictions. At the same
1824 time, ESMFold + ProteinTTT enables more accurate predictions compared to the original ESMFold
1825 (Figure A6 suggests that ESMFold + ProteinTTT significantly improves predictions in almost 40% of
1826 challenging sequences). It means that applying ESMFold + ProteinTTT to these remaining sequences
1827 could significantly expand the metagenomic atlas characterized by ESMFold. Here, we illustrate this
1828 on a similar case study by applying ESMFold + ProteinTTT to more than 300 thousand viral proteins
1829 in BFVD (Section 5.2)

1831 I LIMITATIONS AND FUTURE WORK

1832 We see two main limitations of the current version of ProteinTTT, which we discuss in detail below.

1833 **Extension to other model types and tasks.** The current form of the method is only applicable
1834 to protein language models (PLMs), i.e., Transformer-based (Vaswani, 2017) models pre-trained
1835

1836 using bidirectional masked language modeling (Rives et al., 2021) or autoregressive next-token
1837 prediction (Nijkamp et al., 2023). Nevertheless, the concept of test-time training can also be extended
1838 to many other models in computational biology, which presents exciting opportunities for future
1839 research, as our work demonstrates the high potential of this paradigm for the field of computational
1840 biology. For instance, our central experiments in Section 4.1 use ESMFold (Lin et al., 2023), which
1841 is known to often underperform (Lin et al., 2023) more specialized multiple sequence alignment
1842 (MSA)-based structure predictors such as AlphaFold2 (Jumper et al., 2021), AlphaFold-Multimer
1843 (Evans et al., 2021), AlphaFold3 (Abramson et al., 2024), or Boltz-2 (Passaro et al., 2025).

1844 Nevertheless, all of these models can also be extended with test-time training akin to ProteinTTT.
1845 AlphaFold2, and subsequently AlphaFold-Multimer, use masked modeling of MSA as one of the
1846 training objectives to learn powerful pairwise representations in Evoformer. The Evoformer backbone
1847 could therefore be updated at test time to obtain more powerful representation of one input MSA
1848 at a time using the ProteinTTT objective (Section 3.1). While AlphaFold3 and Boltz-2 do not use
1849 masked modeling, they can still be customized in a self-supervised way, for example using an
1850 optimization through distogram (Cho et al., 2025). Implementing the variants of ProteinTTT for
1851 the aforementioned models could enable customized structure prediction of protein multimers and
1852 protein-ligand complexes.

1853 Beyond structure prediction, test-time customization could also benefit *de novo* protein design. Our
1854 results with autoregressive ProGen2 on fitness prediction suggest that ProteinTTT can improve
1855 sequence design (Table 2). Similarly, although our experiments with ESM3 are currently conducted
1856 in the context of structure prediction (Table 1), ProteinTTT can be straightforwardly applied to
1857 ESM3 for protein design tasks such as inverse folding or structure inpainting by applying ProteinTTT
1858 to the corresponding ESM3 input tracks. Furthermore, BoltzGen can be extended with test-time
1859 customization in a manner analogous to Boltz-2, discussed above, due to their shared architecture.
1860 Performing ProteinMPNN (Dauparas et al., 2022) customization on part of a protein to guide
1861 generation of the remaining structure or its binder, as well as customizing RFdiffusion (Watson et al.,
1862 2023) to a target structure for binder design, represent promising opportunities in protein design with
1863 the potential for higher success rates.

1864 **Better control over failure cases.** The failure modes of ProteinTTT are not yet fully understood.
1865 For instance, combining ESMFold with ProteinTTT decreases performance for several proteins in
1866 the CAMEO test set (Figure A6). A detailed analysis of these cases shows that the degradation
1867 can be attributed to ambiguity in the evaluation itself (Figure A13). These examples illustrate the
1868 challenge of identifying a general reason for the occasional degradation of performance. As discussed
1869 in Appendix H.2, confidence functions (such as pLDDT in structure prediction) allow effectively
1870 eliminating overfitting to a single protein and thereby mitigating such failure cases, making confidence
1871 prediction an essential component of customization.

1872 While confidence functions begin to emerge across tasks, such as fitness prediction (Gurev et al., 2025;
1873 Nijkamp et al., 2023) and inverse folding (Shuai et al., 2025), they are not yet universally available
1874 for use with ProteinTTT. In particular, for fitness (Section 4.2) and function (Section 4.3) prediction,
1875 controlling failure cases remains more challenging due to the absence of a reliable confidence metric.
1876 This motivates the development of general, task-agnostic, unsupervised confidence measures (for
1877 example, perplexity-based estimates (Gurev et al., 2025)) or a dedicated confidence prediction module
1878 within ProteinTTT (Abramson et al., 2024; Jumper et al., 2021). Another promising direction is
1879 deriving confidence estimates from mechanistic interpretability of protein language models (Hübötter
1880 et al., 2025; Simon & Zou, 2025; Zhang et al., 2024).

1881
1882
1883
1884
1885
1886
1887
1888
1889

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

Table A4: **ProteinTTT performance on ProteinGym depending on MSA depth.** MSA depth reflects the number of available proteins similar to the target protein and, when using large protein language models, can be interpreted as a measure of the representation of similar proteins in the training data (Appendix F.2.2). Customization with ProteinTTT primarily improves performance on difficult targets, with low MSA depth. Standard deviations are calculated over 5 random seeds but are omitted in the right panel for brevity, where the maximum standard deviation does not exceed 0.0004.

	Avg. Spearman \uparrow	Spearman by MSA depth \uparrow		
		Low depth	Medium depth	High depth
ESM2 (35M) (Lin et al., 2023)	0.3211	0.2394	0.2707	0.451
ESM2 (35M) + ProteinTTT (Ours)	0.3407 \pm 0.00014	0.2445	0.3144	0.4598
ProGen2-small (151M) (Nijkamp et al., 2023)	0.3255	0.2974	0.3136	0.3765
ProGen2-small (151M) + ProteinTTT (Ours)	0.3591 \pm 0.0002	0.3319	0.3636	0.3917
SaProt (35M) (Su et al., 2023)	0.4062	0.3234	0.3921	0.5057
SaProt (35M) + ProteinTTT (Ours)	0.4106 \pm 0.00004	0.3253	0.3972	0.5091
ESM2 (650M) (Lin et al., 2023)	0.4139	0.3346	0.4063	0.5153
ESM2 (650M) + ProteinTTT (Ours)	0.4153 \pm 0.00003	0.3363	0.4126	0.5075
SaProt (650M) (Su et al., 2023)	0.4569	0.3947	0.4502	0.5448
SaProt (650M) + ProteinTTT (Ours)	0.4583 \pm 0.00001	0.3954	0.4501	0.5439
ProSST (K=2048) (Li et al., 2024)	0.5068	0.4731	0.5107	0.5749
ProSST (K=2048) + ProteinTTT (Ours)	0.5087 \pm 0.00004	0.4809	0.5104	0.5750

Table A5: **Performance of ProteinTTT on the MaveDB dataset.** In this work, we use our newly constructed MaveDB dataset as a validation fold for tuning the ProteinTTT hyper-parameters for fitness prediction. For computational efficiency, we only select a subset of 50 proteins (Appendix F.2.1) and do not run customization across multiple random seeds to estimate standard deviations. The performance shown was calculated by first aggregating correlations per assay, and then per protein (some assays correspond to the same protein).

	Avg. Spearman \uparrow
ESM2 (35M) (Lin et al., 2023)	0.4458
ESM2 (35M) + ProteinTTT (Ours)	0.4593
ESM2 (650M) (Lin et al., 2023)	0.4568
ESM2 (650M) + ProteinTTT (Ours)	0.4604
SaProt (650M) (Su et al., 2023)	0.4926
SaProt (650M) + ProteinTTT (Ours)	0.4926
SaProt (35M) (Su et al., 2023)	0.5251
SaProt (35M) + ProteinTTT (Ours)	0.5271
ProSST (K=2048) (Li et al., 2024)	0.5444
ProSST (K=2048) + ProteinTTT (Ours)	0.5462

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

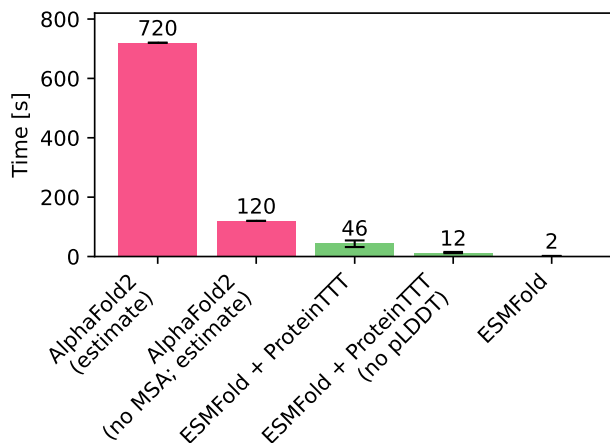


Figure A4: **Running time of ESMFold + ProteinTTT.** For ESMFold and its variants, the median and interquartile ranges of running times on the CAMEO test set are shown using a single NVIDIA A100 GPU. For AlphaFold2, we use estimates from Lin et al. (2023). Specifically, a forward pass through AlphaFold2 is approximately 60 times more computationally expensive than ESMFold (e.g., AlphaFold2 (no MSA; estimate): $2 \times 60 = 120$ seconds), with additional MSA construction taking at least 10 minutes using standard pipelines (AlphaFold2 (estimate): $2 \times 60 + 10 \times 60 = 720$ seconds). ESMFold + ProteinTTT (30 steps) involves LoRA parameter updates, along with forward passes at each customization step to estimate pLDDT and select the structure with the highest predicted confidence. Disabling pLDDT significantly reduces computational overhead (ESMFold + ProteinTTT (no pLDDT) compared to ESMFold + ProteinTTT), but may require careful parameter tuning (Appendix H.2). Overall, ESMFold + ProteinTTT maintains the speed advantage of ESMFold, and is at least an order of magnitude faster than AlphaFold2.

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

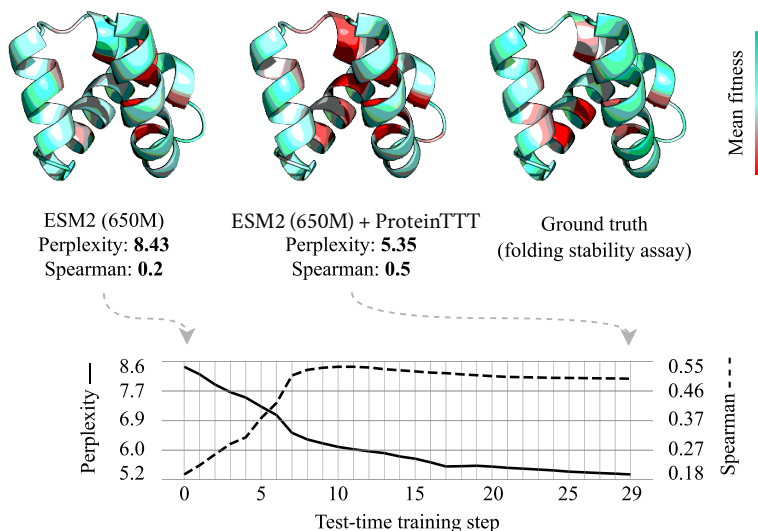


Figure A5: **Example of protein fitness prediction upon single-sequence model customization with ProteinTTT.** Fitness predictions from ESM2 (650M) show poor correlation with experimental fitness values in the ProteinGym test set measured by the stability assay “UBR5_HUMAN_Tsuboyama_2023_1I2T” (Tsuboyama et al., 2023) (left). ESM2 + ProteinTTT achieves significantly higher correlation, likely due to improved detection of disruptive mutations in the protein core that impact protein stability (middle). The ground-truth fitness data aligns with the customized model, showing that residues crucial for stability (i.e., having negative mean fitness) are concentrated in the protein core (right). Residue colors represent the mean fitness upon all single-point substitutions (with the exception of several missing mutations in the ground-truth data), with red indicating residues where mutations have detrimental effects on average.

2052
 2053
 2054
 2055
 2056
 2057
 2058
 2059
 2060
 2061
 2062
 2063
 2064
 2065
 2066
 2067
 2068
 2069
 2070
 2071
 2072
 2073
 2074
 2075
 2076
 2077
 2078
 2079
 2080
 2081
 2082
 2083
 2084
 2085
 2086
 2087
 2088
 2089
 2090
 2091
 2092
 2093
 2094
 2095
 2096
 2097
 2098
 2099
 2100
 2101
 2102
 2103
 2104
 2105

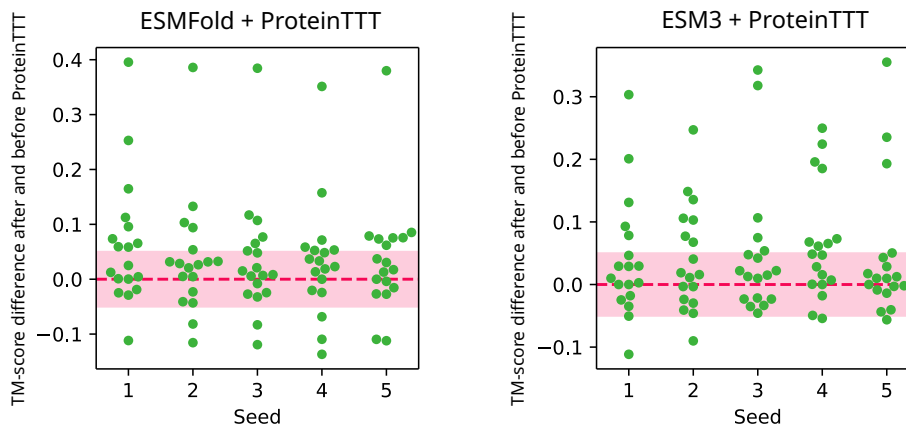


Figure A6: **Per-protein performance of ESMFold + ProteinTTT and ESM3 + ProteinTTT on the CAMEO test set.** The y-axis shows the change in TM-score after applying customization with ProteinTTT, with higher values indicating improvement. The x-axis represents performance across five random seeds. The red dashed line marks no change in TM-score (TM-score difference = 0), and the pink band represents minor changes in TM-score ($-0.05 < \text{TM-score difference} < 0.05$), which we do not consider significant. Each point in the swarm plot corresponds to a single protein from the CAMEO test set. On average, applying ProteinTTT to ESMFold improves the structure predictions for 7 out of 18 proteins, with 2 showing degradation. The rest of the proteins are not significantly affected. Similarly, applying ProteinTTT to ESM3 results in 6 improvements out of 18 proteins, with 1 case of degradation.

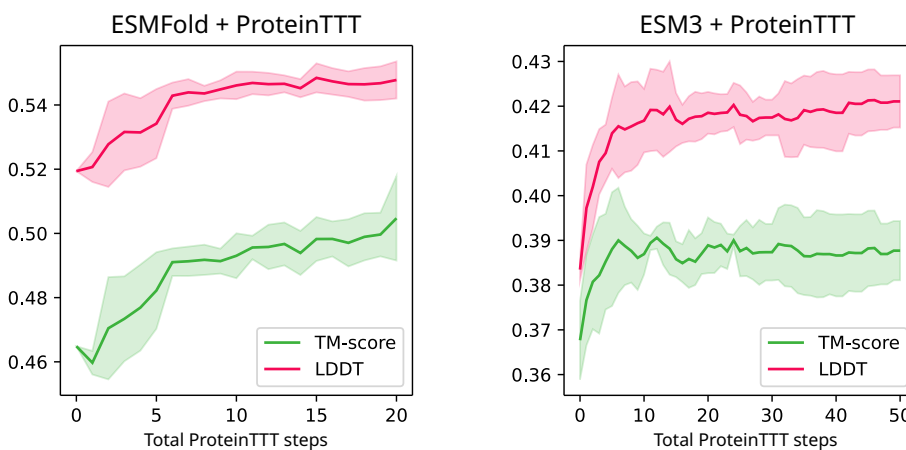


Figure A7: **Test performance of ESMFold + ProteinTTT and ESM3 + ProteinTTT on the CAMEO test set depending on the total number of customization steps.** The x-axis shows the averaged performance across all test proteins, with error bars representing the standard deviation across five random seeds. The y-axis metrics correspond to the structure with the highest pLDDT score up to the given step. While an increased number of ProteinTTT steps generally enhances performance, only a few steps (e.g., five) may suffice to achieve significant performance improvement.

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

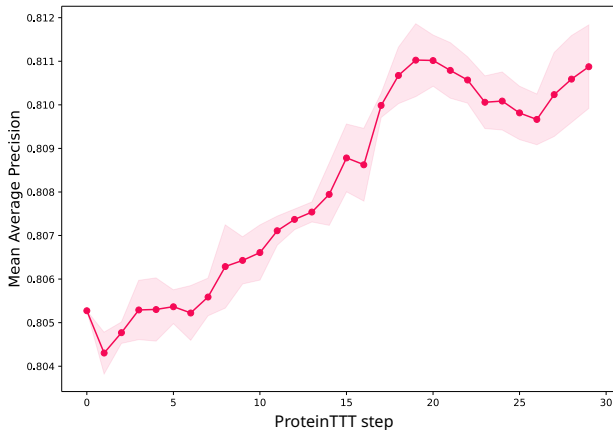


Figure A8: **Test performance of EnzymeExplorer + ProteinTTT across customization steps.** The performance is averaged across all 512 proteins in the dataset, with error bars representing the standard deviation across 5 random seeds.

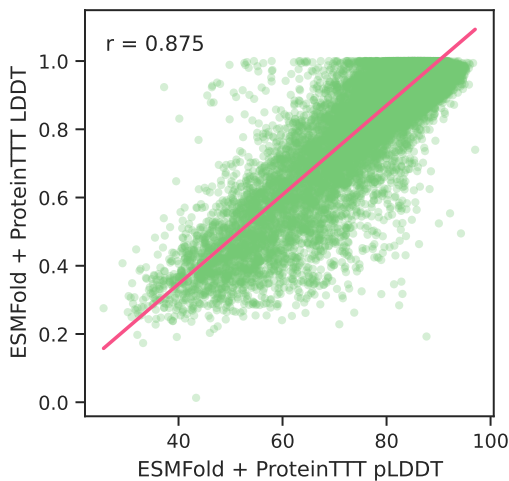


Figure A9: **ESMFold + ProteinTTT pLDDT correlates with ESMFold + ProteinTTT LDDT.** The evaluation was performed on 17,582 AlphaFold2 reference structures from the BFVD database with pLDDT > 90. Here, $r = 0.875$ denotes the Pearson correlation coefficient.

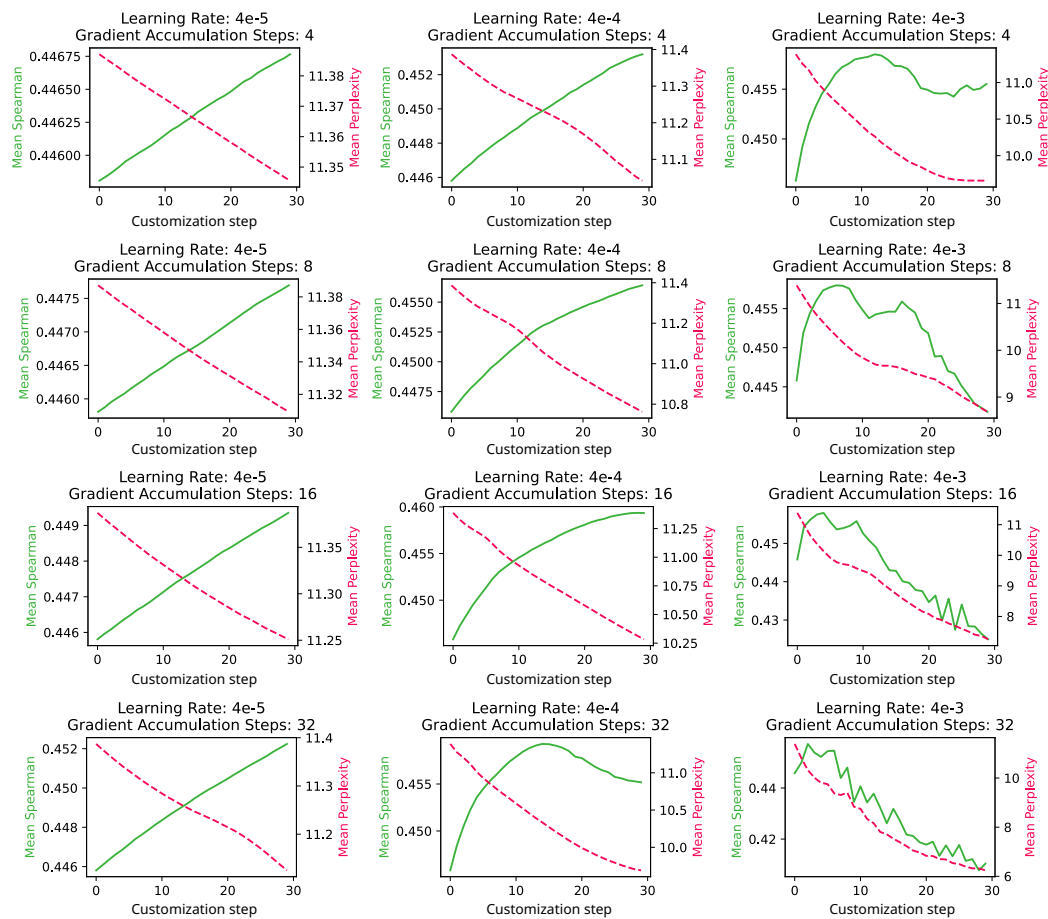


Figure A10: **Dependence on ProteinTTT hyperparameters for customized fitness prediction.** Each plot shows the progression of Spearman correlation (green) increasing alongside a decrease in perplexity (pink) for each customization step, averaged across all assays in the MaveDB validation dataset. The model used is ESM2 (35M) + ProteinTTT, and the grid displays the combinations of different numbers of gradient accumulation steps (i.e., effective batch sizes; shown in rows, increasing from top to bottom) and learning rates (columns, increasing from left to right). As the learning rate increases and the number of gradient accumulation steps grows, the model reaches peak performance more quickly but begins to overfit to a target protein. The optimal hyperparameter combination (learning rate = $4e-4$, gradient accumulation steps = 16) lies near the center of the grid, balancing between underfitting and overfitting to a target protein. Notably, the figure demonstrates that, although ProteinTTT involves three main hyperparameters (batch size, learning rate, and the number of steps), there are effectively only two degrees of freedom controlling the performance of the model. In other words, by keeping the number of steps constant (e.g., 30), the expected performance can be controlled by adjusting the learning rate and the batch size.

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

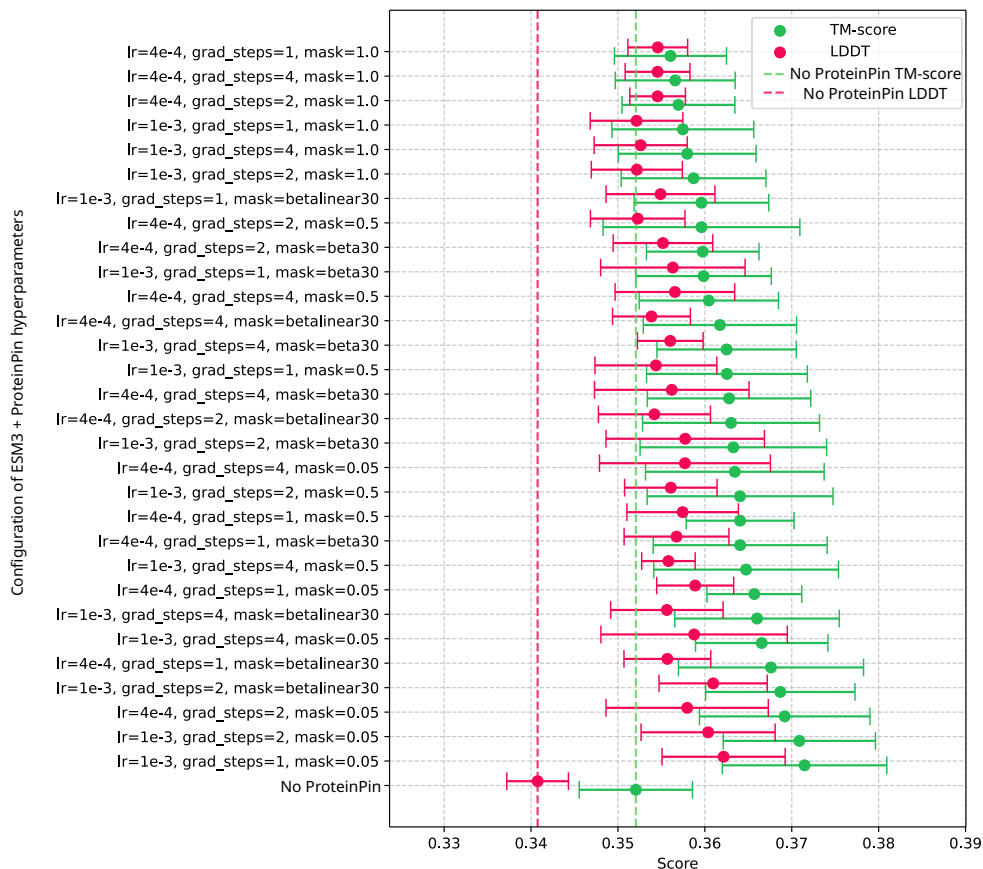


Figure A11: **Hyperparameter search for protein structure prediction with ESM3 + ProteinTTT.**

We conducted a comprehensive grid search based on three key hyperparameters: learning rate (denoted as “lr”), number of gradient accumulation steps (denoted as “grad_steps”; with the batch size of two), and masking strategy (denoted as “mask”). We explored two learning rates, 4e-4 and 1e-3, three gradient accumulation step values of 1, 4, and 16, and five different masking strategies: uniform sampling of 0.05, 0.5, and 1.0 fractions of amino acids, as well as the “beta30” and “betalinear30” distributions proposed in the ESM3 paper (Hayes et al., 2024). Each row in the table presents the mean TM-score and LDDT metrics with standard deviations across five random seeds on the CAMEO validation fold. The last row, denoted as “No ProteinTTT”, shows the performance of ESM3 without customization. The results indicate that ESM3 + ProteinTTT is robust to the choice of hyperparameters and consistently outperforms the base model across all configurations. We selected the configuration from the last row (excluding “No ProteinTTT”) to compute the results on the test fold. For the hyperparameter search, we used 30 customization steps instead of 50 to reduce computation time.

2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321

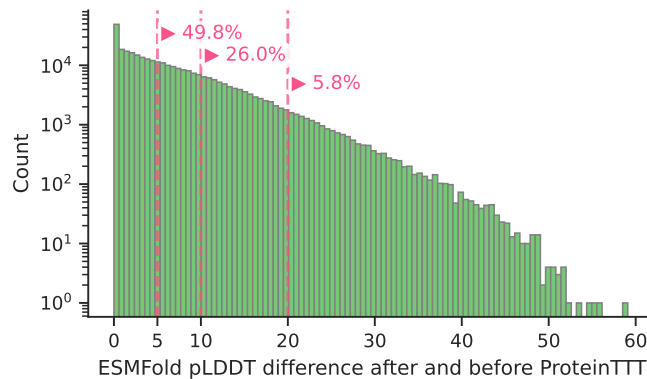


Figure A12: **Magnitude of ESMFold pLDDT improvements after customization with ProteinTTT.** The evaluation is performed on 317,882 proteins from the Big Fantastic Virus Database (BFVD). Percentage annotations indicate the fraction of proteins whose pLDDT increases by at least the corresponding value (e.g., 49.8% of proteins show an improvement of at least 5 pLDDT points).

2322
 2323
 2324
 2325
 2326
 2327
 2328
 2329
 2330
 2331
 2332
 2333
 2334
 2335
 2336
 2337
 2338
 2339
 2340
 2341
 2342
 2343
 2344
 2345
 2346
 2347
 2348
 2349
 2350
 2351
 2352
 2353
 2354
 2355
 2356
 2357
 2358
 2359
 2360
 2361
 2362
 2363
 2364
 2365
 2366
 2367
 2368
 2369
 2370
 2371
 2372
 2373
 2374
 2375

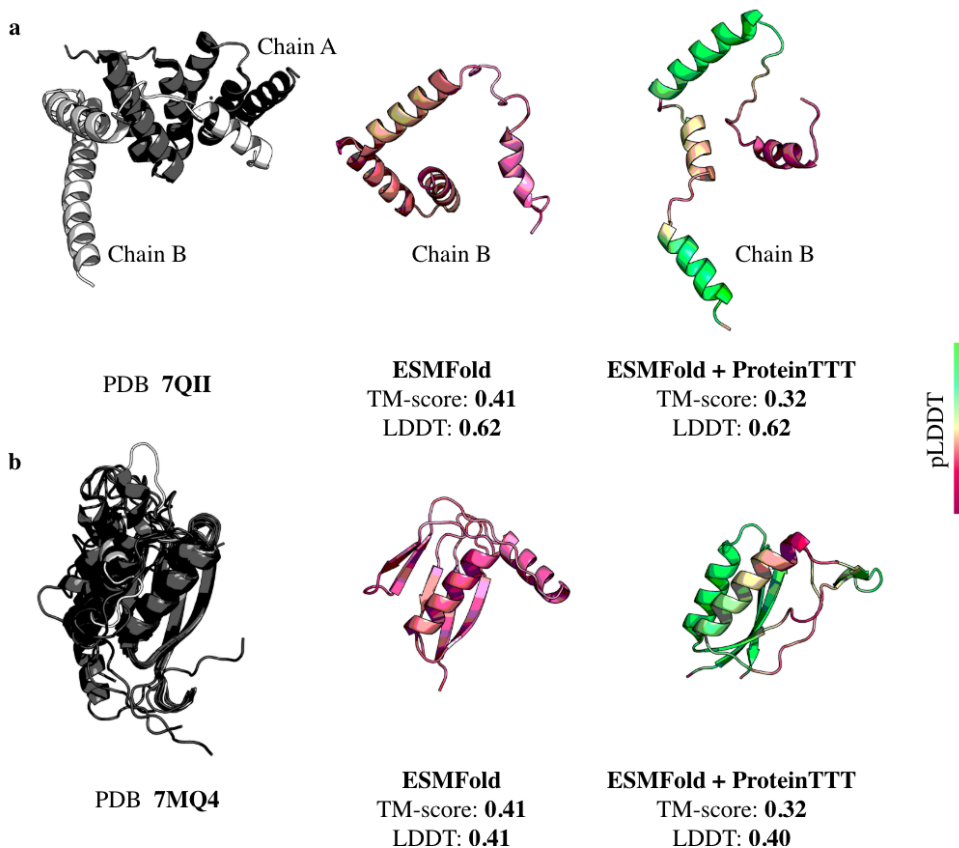


Figure A13: **Detailed analysis of ProteinTTT failure cases on the CAMEO test set.** The figure shows the two entries that consistently exhibit a decrease in TM-score after customization with ProteinTTT across most random seeds (see Figure A6). **(a)** For chain B of PDB entry 7QII (white), the ground-truth structure is part of a dimer in which the conformation of chain B depends on interactions with chain A (black). In the monomeric prediction setting, this context is absent, making the precise helix arrangement inherently ambiguous. Both ESMFold and ESMFold + ProteinTTT correctly capture the helical composition but differ in the global configuration, leading to different TM-scores. **(b)** For chain A of PDB entry 7MQ4 (white), the reference structure is an NMR ensemble with substantial conformational variability (black). Both ESMFold and ESMFold + ProteinTTT recover the stable substructure (right part of the structure in black consisting of a helix surrounded by beta strands), yet produce different conformations in the flexible regions, where multiple arrangements are plausible.