

TRAINING ON TEST PROTEINS IMPROVES FITNESS, STRUCTURE, AND FUNCTION PREDICTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Data scarcity and distribution shifts often hinder the ability of machine learning models to generalize when applied to proteins and other biological data. Self-supervised pre-training on large datasets is a common method to enhance generalization. However, striving to perform well on all possible proteins can limit model’s capacity to excel on any specific one, even though practitioners are often most interested in accurate predictions for the individual protein they study. To address this limitation, we propose an orthogonal approach to achieve generalization. Building on the prevalence of self-supervised pre-training, we introduce a method for self-supervised fine-tuning at test time, allowing models to adapt to the test protein of interest on the fly and without requiring any additional data. We study our test-time training (TTT) method through the lens of perplexity minimization and show that it consistently enhances generalization across different models, their scales, and datasets. Notably, our method leads to new state-of-the-art results on the standard benchmark for protein fitness prediction, improves protein structure prediction for challenging targets, and enhances function prediction accuracy.

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 demonstrated 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; Yang et al., 2024; Kouba et al., 2023). However, biologists often focus their research on individual proteins or protein complexes involved for example in metabolic disorders (Ashcroft et al., 2023; Gunn & Neher, 2023), oncogenic signalling (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.

Nonetheless, general machine learning models for proteins often struggle to generalize to individual case studies due to data scarcity (Bushuiev et al., 2023; Chen & Gong, 2022) and distribution shifts (Tagasovska et al., 2024; Feng et al., 2024). Bridging the gap between broad, dataset-wide optimizations and the precision required for studying single proteins in

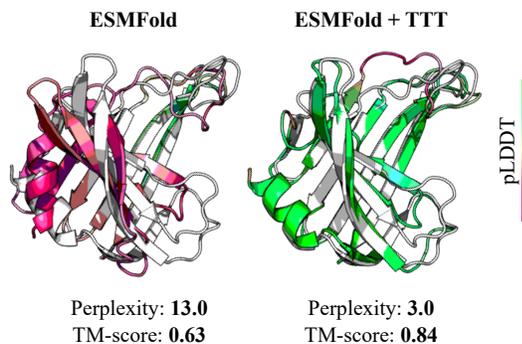


Figure 1: Example of test-time training (TTT) applied to protein folding. ESMFold poorly predicts the structure of the CASP14 target T1074 (shown in white) because the underlying language model ESM2 poorly fits the sequence, as indicated by the high perplexity (Fig. 2E in Lin et al. (2023) and the left panel here). Self-supervised test-time training of ESM2 on the single sequence of T1074 minimizes the perplexity, leading to improved structure prediction (better TM-score alignment and higher pLDDT predicted confidence). The same test-time training approach is also broadly applicable to other tasks, such as protein fitness and function prediction.

054 practical applications remains a critical challenge in integrating machine learning into biological
055 research (Sapoval et al., 2022).

056 By contrast, in other application domains of machine learning, such as computer vision and natural
057 language processing, customization and adaptation approaches have emerged as powerful tools to
058 improve model performance in specific contexts (Ruiz et al., 2023; Hardt & Sun, 2023). Drawing
059 inspiration from the test-time training (TTT) approach developed in computer vision to mitigate
060 distribution shifts (Sun et al., 2020; Gandelsman et al., 2022), in this work we propose the TTT
061 approach for proteins. Our method enables adapting protein models to one protein at a time, on the
062 fly, and without the need for additional data. Given a model that has been pre-trained using masked
063 language modeling, our method minimizes the perplexity of the model on a given test protein through
064 self-supervised fine-tuning, which, in turn, results in improved downstream performance without
065 updating the downstream task head.

066 The prevalence of masked modeling in protein machine learning makes our method broadly applicable
067 to various downstream tasks. Empirically, we demonstrate its effectiveness across three key challenges
068 in protein machine learning. First, TTT achieves state-of-the-art results on the ProteinGym dataset
069 (Notin et al., 2024), a well-established benchmark for protein fitness prediction. Second, TTT
070 enhances protein structure predictions with ESMFold (Lin et al., 2023) and ESM3 (Hayes et al.,
071 2024) on challenging targets. Third, the application of TTT to protein function predictors results in
072 improved classification of terpene synthase (TPS) substrates and protein subcellular localization.

073 In summary, the key contributions of this work are three-fold:

- 074 1. Motivated by the generalization challenges and distribution shifts prevalent in protein
075 machine learning, we introduce a new test-time training (TTT) method¹ that enables models
076 to adapt to individual proteins on the fly and without requiring additional data.
- 077 2. We establish a link between our TTT approach and perplexity minimization, providing an
078 insight into why this approach enhances model effectiveness.
- 079 3. We empirically validate TTT, achieving state-of-the-art results in protein fitness prediction,
080 improving the protein structure prediction capabilities of well-established folding mod-
081 els, and enhancing protein function predictions in the tasks of terpene synthase substrate
082 classification and protein localization prediction.

083 2 BACKGROUND AND RELATED WORK

084 In this section, we present the context and related work that highlight the [rationale](#), feasibility, and
085 broad applicability of test-time training (TTT) in the domain of machine learning on proteins. The
086 widespread adoption of Y-shaped architectures relying on masked modeling enables the development
087 of a general method for adapting protein models at test time via masking-based self-supervised
088 fine-tuning.

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

099 The backbone f is typically a large neural network pre-trained in a self-supervised way on a large
100 dataset using a smaller pre-training projection head g (Hayes et al., 2024). The fine-tuning head h ,
101 however, depends on the application. In some cases, h is a large neural network, repurposing the
102 pre-trained model entirely (Watson et al., 2023; Lin et al., 2023); in others, h is a minimal projection

103 ¹<https://github.com/anton-bushuiev/ProteinTTT>

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)). In some cases, the fine-tuning head h may also be a machine learning algorithm other than a neural network (Samusevich et al., 2024).

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, most methods employ various forms of masked modeling, i.e., optimizing the model weights to accurately reconstruct missing parts of proteins, regardless of the downstream application. Masked modeling pre-training underpins models for protein structure (Lin et al., 2023) and function (Samusevich et al., 2024) prediction, as well as for protein design (Hayes et al., 2024). [For example, in AlphaFold2, a significant part of the loss function weight is put onto masked modeling of multiple sequence alignments \(MSAs\)](#) (Jumper et al., 2021), and the model has been effectively fine-tuned for various tasks beyond structure prediction (Jing et al., 2024; Cheng et al., 2023; Motmaen et al., 2023).

Masked modeling is a dominant pre-training objective not only across different tasks but also across various protein representations. Sequence models applied to proteins are typically pre-trained to predict randomly masked amino acids in a random or autoregressive manner (Lin et al., 2023; Rao et al., 2021; Elnaggar et al., 2023; Madani et al., 2023; Ferruz et al., 2022; Rives et al., 2021; Rao et al., 2020). Models utilizing graph neural networks or 3D convolutions on protein structures are also commonly pre-trained to fill in missing structural fragments (Dieckhaus et al., 2024; Diaz et al., 2023; Bushuiev et al., 2023; Hsu et al., 2022; Shroff et al., 2020). The most recent approaches combine both sequential and structural information under masked modeling (Hayes et al., 2024; Su et al., 2023; Heinzinger et al., 2023).

Model adaptation. In many scenarios, machine learning models for proteins benefit from being adapted to a specific protein of interest. This adaptation is commonly achieved in two ways: either via additional input features or via protein-specific fine-tuning. Multiple sequence alignments (MSAs) containing sequences similar to the target protein provide a common way of supplying a model with protein-specific features (Abramson et al., 2024; Jumper et al., 2021; Rao et al., 2021). Another approach for injecting protein-specific knowledge into the model is standard supervised fine-tuning (i.e., via the $h \circ f$ track) on protein-specific data (Notin et al., 2024; Kirjner et al., 2023; Rao et al., 2019). An alternative is self-supervised fine-tuning (i.e., via the $g \circ f$ track) on proteins from the MSA (Notin et al., 2022b; Frazer et al., 2021; Alley et al., 2019) or on proteins sharing another property with the target protein, such as common family (Sevgen et al., 2023) or class (Samusevich et al., 2024). However, constructing MSAs is time-consuming (Fang et al., 2023), and similar proteins may not be available for many targets (Durairaj et al., 2023; Lin et al., 2023).

Here, we propose an extreme case of self-supervised fine-tuning: learning from a single target protein, without the need for any additional data. To the best of our knowledge, this approach has not been employed in the field of machine learning applied to biology; however, similar methods have been developed in computer vision (Chi et al., 2024; Wang et al., 2023; Xiao et al., 2022; Karani et al., 2021) and natural language processing (Hardt & Sun, 2023; Ben-David et al., 2022; Banerjee et al., 2021). The paradigm of test-time training (TTT), developed to mitigate distribution shifts in computer vision applications (Gandelsman et al., 2022; Sun et al., 2020), is a main inspiration for our work. Here, we demonstrate that TTT is highly relevant for machine learning on proteins even without the presence of explicit distribution shift. We investigate the link of TTT to perplexity minimization and show that TTT improves performance on several downstream tasks.

3 TEST-TIME TRAINING (TTT) ON PROTEINS

As discussed in the previous section, many machine learning models for proteins employ Y-shaped architectures, consisting of a backbone f with a self-supervised head g and a supervised head h . This design facilitates the use of self-supervised fine-tuning across various tasks and models. Notably, most of these models leverage masked modeling as a pre-training objective, which enables the introduction of a broadly applicable test-time training (TTT) method based on masking. Our method adapts models to specific test proteins through masked modeling (Figure 2). In this section, we first formally define the proposed TTT approach (Section 3.1), followed by its application to a range of well-established models (Section 3.2). Finally, we provide an insight into the effectiveness of our method by linking it to perplexity minimization (Section 3.3).

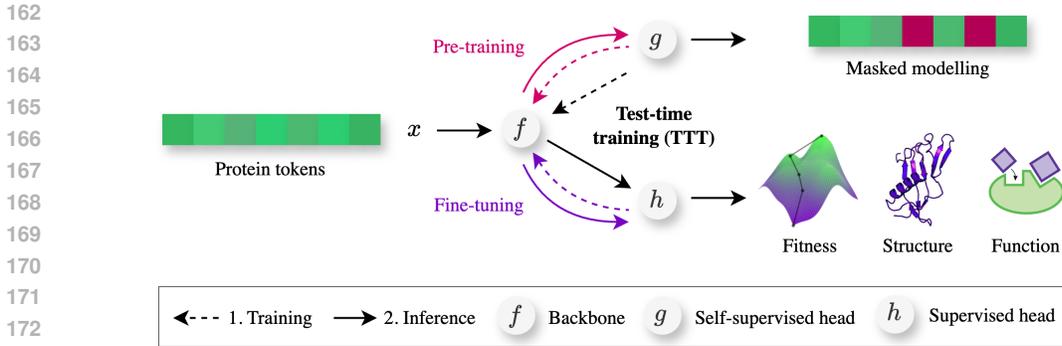


Figure 2: **Overview of our test-time training (TTT) for proteins.** Test-time training for proteins builds on the prevalence of Y-shaped architectures relying on masked modeling (i.e., self-supervised masking-based pre-training of $g \circ f$ followed by supervised fine-tuning of $h \circ f$, sharing the backbone f). Given a single test protein x , TTT adapts the backbone f to the protein using self-supervised fine-tuning. This adaptation leads to better generalization for the downstream task, such as protein fitness, structure, or function prediction.

3.1 SELF-SUPERVISED FINE-TUNING ON TEST PROTEINS

At test time, we assume a Y-shaped model with a backbone f that has been pre-trained via the self-supervised track $g \circ f$, followed by task-specific fine-tuning through the supervised track $h \circ f$. The goal of test-time training (TTT) is to adapt the backbone f to a single test example x before performing test-time inference on a downstream task via the supervised track.

To achieve this, we first fine-tune **all layers of** the backbone f using the self-supervised track $g \circ f$ on the single example x . This step customizes the backbone f to the test sample x , and, as demonstrated in Section 4, enhances the generalization of $h \circ f$ without modifying the weights of the task-specific head h . Figure 2 illustrates our method. Although the concept of TTT is relatively simple, it involves several important design choices, such as selecting the optimizer and efficiently fine-tuning large backbones, which we describe in the following paragraphs.

Training objective. We fine-tune $g \circ f$ on a test sample x via minimizing the masked language modeling objective (Devlin, 2018; Rives et al., 2021):

$$\mathcal{L}(x) = \mathbb{E}_M \left[\sum_{i \in M} -\log p(x_i | x_{\setminus M}) \right], \quad (1)$$

where x denotes a sequence of protein tokens (typically amino acid types), and \mathbb{E}_M represents the expectation over randomly sampled masking positions M . The loss function $\mathcal{L}(x)$ maximizes the log-probabilities $\log p(x_i | x_{\setminus M})$ of the true tokens x_i at the masked positions $i \in M$ in the partially masked sequence $x_{\setminus M}$. Please note that here we focus on bi-directional masked modeling models, which employ random masking, but the method can be straightforwardly extended to models employing autoregressive masking.

In practice, \mathbb{E}_M can 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 test-time training, we replicate the masking distribution used during the pre-training. If relevant, we also replicate other pre-training tricks, 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. We minimize the loss defined in Equation (1) using stochastic gradient descent (SGD) with zero momentum and zero weight decay (Ruder, 2016). While a more straightforward option might be to use the optimizer state from the final pre-training step, this approach is often impractical because the optimizer parameters are usually not provided with the pre-trained model

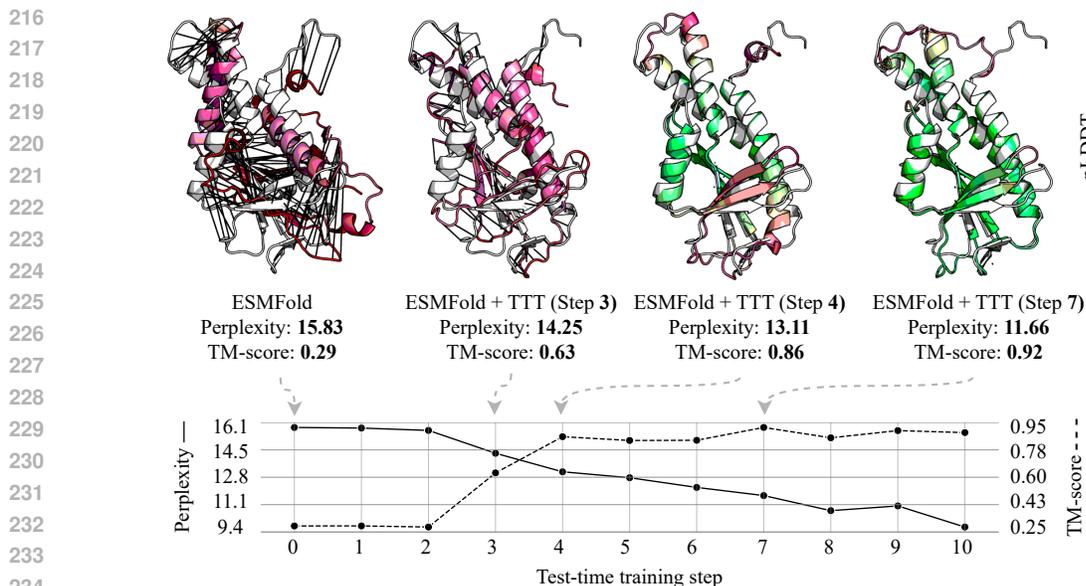


Figure 3: **Test-time training (TTT) improves protein structure prediction by reducing protein sequence perplexity.** ESMFold fails to predict the structure of chain B from PDB entry 7EBL in the CAMEO validation set, as shown at TTT step 0, where the perplexity is high and the TM-score is low. By applying TTT on the single target sequence, the model iteratively improves the structure prediction quality, as demonstrated by the increasing TM-score, associated with reduced perplexity. At 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 ESMFold + TTT.

weights (Hayes et al., 2024; Lin et al., 2023). Moreover, many models are pre-trained using the Adam optimizer (Kingma & Ba, 2015) or its variants (Loshchilov & Hutter, 2019). However, it has been shown that Adam results in less predictable behavior of test-time training (TTT) compared to the SGD optimizer, possibly due to its more exploratory behavior (Gandelsman et al., 2022).

Because each TTT experiment assumes only one test example available, we are not able to halt the training using early stopping on any validation sample. Therefore, for each choice of task-specific f and h , we tune the optimal number of TTT steps using the entire validation set beforehand or rely on available performance estimates (e.g., pLDDT in the case of protein structure prediction; Section 4.2) to select the optimal number of optimization steps.

Fine-tuning large models. We aim for test-time training to be applicable on the fly, i.e., without the need for any pre-computation and on a single GPU with a minimum computational overhead. Since state-of-the-art models for many protein-oriented tasks are typically large, with up to billions of parameters, our aim presents two key challenges. First, when using pre-trained transformers on a single GPU, even for the forward pass, the batch size is typically limited to only several samples due to the quadratic complexity of the inference (Vaswani, 2017). Second, for the backward pass, even a batch size of one is not always feasible for large models. To address the first challenge, we perform forward and backward passes through a small number of training examples and accumulate gradients to simulate updates with any batch size. We address the second challenge by employing low-rank adaptation (LoRA, Hu et al. (2021)), which in practice enables fine-tuning of any model for which a forward pass on a single sample is feasible, due to a low number of trainable parameters.

3.2 INFERENCE ON DOWNSTREAM TASKS

Once the backbone f is adapted to a test protein via self-supervised fine-tuning, it can be used in conjunction with a pre-trained downstream head h , as $h \circ f$. The key idea of TTT is not to update the head h during test time, but rather to leverage improved input representations from f .

Since Y-shaped architectures are prevalent in protein machine learning, TTT can be straightforwardly applied to numerous tasks in protein research. In this work, we address three primary challenges: protein fitness, structure, and function prediction, applying TTT to corresponding well-established models. For fitness prediction, we apply TTT to ESM2 (Lin et al., 2023) and SaProt (Su et al., 2023); for folding, we apply it to ESMFold (Lin et al., 2023) and ESM3 (Hayes et al., 2024); and for function prediction, we apply TTT to ESM-1v-based (Meier et al., 2021) TerpeneMiner (Samusevich et al., 2024) and ESM-1b-based (Rives et al., 2021) Light attention (Stärk et al., 2021).

In all the models we consider, f is a transformer encoder that takes a protein sequence as input (except for SaProt, which also uses structural tokens), while g is a masked language modeling head (a layer mapping token embeddings to amino acid types). The downstream task heads h , however, vary significantly across tasks. For fitness prediction, h outputs a single value for a mutated sequence, measuring how well the protein supports an organism’s functioning. **Both ESM2 and SaProt perform zero-shot inference using $h \circ f$ via log odds from g , with h functioning as a simple adaptation of g without introducing additional parameters.** For structure prediction, h is a protein structure decoder: in ESMFold, it is an AlphaFold2-like structure prediction module (Jumper et al., 2021), while in ESM3, it is a VQ-VAE decoder (Razavi et al., 2019). The function predictors are classification models: in TerpeneMiner (Samusevich et al., 2024), h is a random forest that outputs substrate probabilities, and in Light attention (Stärk et al., 2021), h is a light attention module predicting localization class probabilities. Detailed descriptions of the models and their TTT adaptation are provided in Appendix A.

3.3 JUSTIFICATION FOR TEST-TIME TRAINING VIA PERPLEXITY MINIMIZATION

While the approach of test-time training has been extensively investigated in computer vision and other domains, the reasons behind its effectiveness remain unclear (Liu et al., 2021; Zhao et al., 2023). **Here, we offer a potential justification for the effectiveness of TTT by linking it to perplexity minimization within the context of protein sequence modeling.**

Perplexity has traditionally been used in natural language processing to evaluate how well models comprehend test 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:

$$\text{Perplexity}(x) = \exp\left(\frac{1}{|x|} \sum_{i=1}^{|x|} -\log p(x_i|x_{\setminus i})\right), \quad (2)$$

where $|x|$ is the length of the input protein sequence x and $p(x_i|x_{\setminus i})$ 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 understands, on average, positions within a given sequence. A perplexity value of 1 indicates that the model perfectly understands 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 provide an example, we analyze the correlation between perplexity and structure prediction performance (Figure 4). A strong correlation suggests that

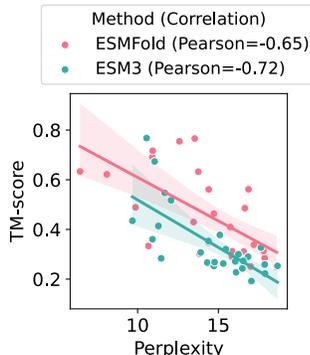


Figure 4: The quality of protein structure prediction, as measured by TM-score, correlates with perplexity of the underlying language model on the challenging targets from the CAMEO validation set. Higher TM-scores are associated with lower perplexity, indicating that better predictions are linked to 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).

324 reducing a model’s perplexity on a single test sample x can lead to improved performance on the
325 downstream task (Figure 3; Figure 12).

326 Since we consider only a single test example x , the minimization of the masked language mod-
327 eling loss $\mathcal{L}(x)$ (Equation (1)) on this example is directly linked to minimizing the perplexity
328 $\text{Perplexity}(x)$ (Equation (2)). For instance, in the case of a single masked position (i.e., $|M| = 1$),
329 the loss is equal to the logarithm of perplexity. More generally, it can be shown formally that by
330 minimizing the masked language modeling objective, one learns to approximate the conditional
331 marginals of the language (of proteins), including the leave-one-out probabilities evaluated in perplex-
332 ity (Hennigen & Kim, 2023). As a result, applying test-time training (TTT) through $g \circ f$ enhances
333 the representation of the test protein in the backbone f , leading to improved downstream performance
334 via the fine-tuning track $h \circ f$.

336 4 EXPERIMENTS

337 Building on the broad applicability of our test-time training (TTT) approach, we apply it to three
338 example downstream tasks in protein machine learning: fitness, structure, and function prediction.
339 The experimental setup and results for each task are presented in the following subsections.

342 4.1 PROTEIN FITNESS PREDICTION

343 Protein fitness refers to the ability of a protein to efficiently perform its biological function, which
344 is determined by its structure, stability, and interactions with other molecules. Predicting protein
345 fitness allows researchers to understand how mutations affect protein function, aiding in protein
346 engineering (Notin et al., 2024). In this paper, we demonstrate that applying test-time training (TTT)
347 to representative models, such as ESM2 (Lin et al., 2023) and SaProt (Su et al., 2023), enhances
348 their protein fitness prediction capabilities. ESM2 is a protein language model trained on protein
349 sequences, while SaProt is an extension of ESM2 that incorporates 3D information via additional
350 structural tokens encoding structures predicted by AlphaFold2 (Jumper et al., 2021).

352 **Evaluation Setup.** We evaluate the models using ProteinGym, [state-of-the-art](#) benchmark for
353 fitness prediction (Notin et al., 2024), focusing specifically on its [well-established](#) zero-shot variant.
354 [The zero-shot nature of this benchmark enables us to validate TTT in a simplified setting with a](#)
355 [minimalist head \$h\$, which is complementary to the other tasks described below.](#) Since the zero-shot
356 setup only provides a test set without any data split, we aim to validate TTT on independent data.
357 To achieve this, we create a new fitness prediction dataset mined from MaveDB, a public repository
358 containing datasets from Multiplexed Assays of Variant Effect (MAVEs) (Esposito et al., 2019). The
359 quality of the new dataset is validated by confirming that both ESM2 and SaProt generalize well to
360 the new data, achieving comparable performance (Appendix A).

361 Given a protein and its variants, fitness prediction models output one real value per variant to estimate
362 fitness. ProteinGym uses Spearman correlation between predicted and experimentally measured
363 fitness values as the main evaluation metric for assessing the capabilities of models to score mutations.
364 The correlation is first calculated for each protein and then aggregated per types of measured fitness:
365 activity, binding, expression, organismal fitness, and stability. The final Spearman correlation metric
366 is obtained by averaging across these five categories. We adopt this metric in our benchmarking.

367 In our evaluation, we also include other top-performing baselines on the ProteinGym benchmark:
368 TranceptEVE (Notin et al., 2022b) and GEMME (Laine et al., 2019). TranceptEVE combines
369 language model Tranception (Notin et al., 2022a) with the protein-specific variational autoencoder,
370 EVE, capturing the evolutionary information via MSAs (Frazer et al., 2021). GEMME is a statistical
371 method deriving fitness predictions from evolutionary trees.

372 **Results.** Test-time training (TTT) consistently enhances the protein fitness prediction performance
373 of both ESM2 and SaProt models across varying model scales (35M and 650M parameters) and
374 both datasets, test ProteinGym (Table 1 left) and validation MaveDB (Table 6 in Appendix B.2).
375 Notably, SaProt (650M) + TTT sets a new state-of-the-art on the ProteinGym benchmark, achieving
376 a 40% higher improvement compared to the previous leaderboard update (SaProt (650M) against
377 TranceptEVE L). When examining performance across different phenotype categories, TTT yields

Table 1: **Test-time training (TTT) 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 and TTT is applied to each protein individually. Standard deviations are calculated over 5 random seeds and, for brevity, omitted in the right panel, where the maximum standard deviation does not exceed 0.0004. Methods marked with an asterisk (“*”) are the other top-5 methods in ProteinGym, and the metrics are reproduced from the leaderboard (<https://proteingym.org/benchmarks>).

	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) + TTT (Ours)	0.3407 \pm 0.00014	0.3407	0.2942	0.3550	0.2403	0.4733
SaProt (35M) (Su et al., 2023)	0.4062	0.3721	0.3568	0.4390	0.2879	0.5749
SaProt (35M) + TTT (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) + TTT (Ours)	0.4153 \pm 0.00003	0.4323	0.3376	0.4168	0.3702	0.5195
TranceptEVE S* (Notin et al., 2022b)	0.4519	0.4750	0.3957	0.4426	0.4491	0.4973
GEMME* (Laine et al., 2019)	0.4547	0.4820	0.3827	0.4382	0.4517	0.5187
TranceptEVE M* (Notin et al., 2022b)	0.4548	0.4792	0.3858	0.4525	0.4538	0.5025
TranceptEVE L* (Notin et al., 2022b)	0.4559	0.4866	0.3758	0.4574	0.4597	0.5003
SaProt (650M) (Su et al., 2023)	0.4569	0.4584	0.3785	0.4884	0.3670	0.5919
SaProt (650M) + TTT (Ours)	0.4583 \pm 0.00001	0.4593	0.3790	0.4883	0.3754	0.5896

improvements specifically in the categories where the baseline performance is weakest: “Organismal Fitness”, “Binding”, and “Activity” (Table 1 right). This improvement indicates the ability of TTT to enhance predictions on challenging targets. Additionally, we observe an inverse correlation between the degree of TTT enhancement and the depth of the MSA (i.e., the number of available homologous sequences) available for each test protein, suggesting that TTT primarily improves predictions for proteins with fewer similar sequences available in the training data (Table 5 in Appendix B.1). Interestingly, TTT more effectively enhances the performance of smaller ESM2 and SaProt models compared to their larger variants (Table 1 and Table 6 in Appendix B) and does not require the application of LoRA even for the larger models (Table 4).

4.2 PROTEIN STRUCTURE PREDICTION

Protein structure prediction, also known as protein folding, is the task of predicting 3D coordinates of protein atoms given the amino acid sequence. Arguably, one of the most remarkable applications of machine learning in the life sciences has been in protein folding (Jumper et al., 2021; Lin et al., 2023; Abramson et al., 2024), paving the way for numerous advances in the understanding of biology (Yang et al., 2023; Akdel et al., 2022; Barrio-Hernandez et al., 2023). However, even state-of-the-art protein folding methods struggle to generalize to entirely novel proteins (Kryshtafovych et al., 2023). In this work, we focus on the ESMFold (Lin et al., 2023) and ESM3 (Hayes et al., 2024) models, demonstrating how their performance on challenging targets can be boosted by utilizing TTT.

Evaluation setup. To evaluate the performance of TTT, we use CAMEO, a standard benchmark for protein folding. We use the validation and test folds from Lin et al. (2023), focusing only on challenging targets by filtering them according to standard measures of prediction confidence based on pLDDT and perplexity (Appendix A.2).

Given a protein sequence, the goal of protein folding is to predict 3D coordinates of the protein atoms. To assess the quality of the predicted protein structures with respect to the ground truth structures, we use two standard metrics: TM-score (Zhang & Skolnick, 2004) and LDDT (Mariani et al., 2013). TM-score measures the quality of the global 3D alignment of the target and predicted protein structures, while LDDT is an alignment-free method based on local distance difference tests.

As baseline methods, we use techniques alternative to TTT for improving the performance of the pre-trained base models. In particular, the ESMFold paper proposes randomly masking 15% of

Table 2: **Test-time training (TTT) improves protein structure prediction.** The metrics are averaged across the 18 challenging targets (TTT is applied to each protein individually) 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.

	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 + TTT (Ours)	0.3954 \pm 0.0067	0.4214 \pm 0.0054
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 + TTT (Ours)	0.5047 \pm 0.0132	0.5478 \pm 0.0058

amino acids in a protein sequence, 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 TTT 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. Test-time training (TTT) consistently improves the performance of both the ESMFold and ESM3 models, outperforming the masked prediction (ESMFold + MP) and chain-of-thought (ESM3 + CoT) baselines, as shown in Table 2. Of the 18 most challenging CAMEO test proteins, ESMFold and ESM3 significantly improved the prediction of 7 and 6 structures, respectively, while only slightly disrupting the prediction of 2 and 1 structures, respectively (Figure 9 in Appendix B.1). Most notably, TTT enables accurate structure prediction for targets that are poorly predicted with original base models. For instance, Figure 1 presents a strongly improved structure predicted using ESMFold + TTT 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 TTT refined the structure prediction from a low-quality prediction (TM-score = 0.29) to a nearly perfectly folded protein (TM-score = 0.92). Figure 8 in Appendix B shows that ESMFold + TTT maintains computational efficiency comparable to ESMFold while being orders of magnitude faster than AlphaFold2. Figure 13 in Appendix B additionally demonstrates the robustness of ESM3 + TTT to the choice of hyperparameters.

4.3 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., 2024). While improved structure prediction with TTT (Section 4.2) can already enhance function prediction (Song et al., 2024), we also evaluate TTT 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 terpenoids, the largest class of natural products (Christianson, 2017; Samusevich et al., 2024). Using TTT with TerpeneMiner (Samusevich et al., 2024) 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. (2024) and repurpose 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 a 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 TerpeneMiner paper (Samusevich et al., 2024): 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.

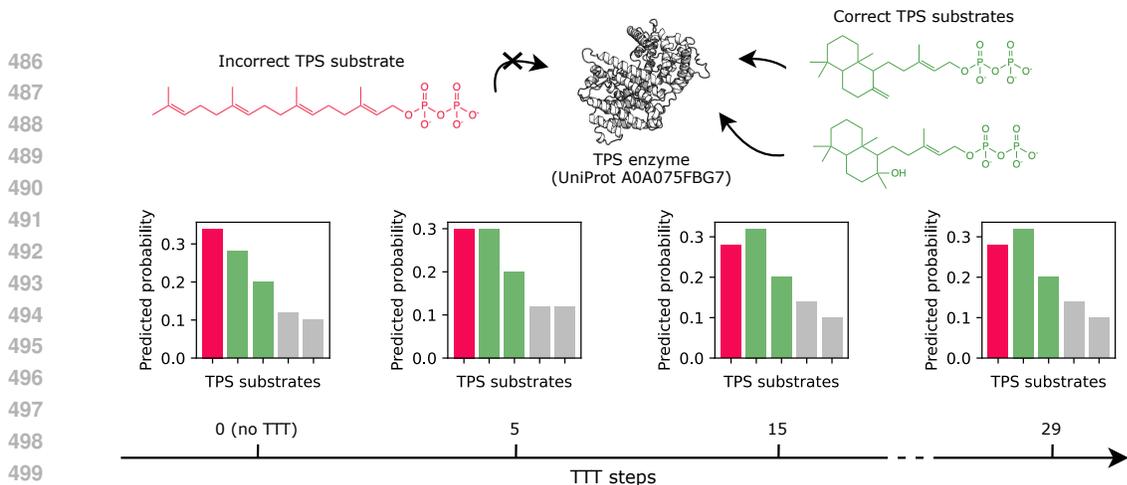


Figure 5: **Test-time training (TTT) enables the correct substrate classification for a terpene synthase (TPS) enzyme.** With progressive test-time training steps of TerpeneMiner + TTT, the probability of the initially misclassified substrate (red) decreases, while the probability of the true substrates (green) increases. The bar plots also display the predicted probabilities for other substrates with non-zero values (grey).

Table 3: **Test-time training (TTT) improves protein function prediction.** For the terpene synthase (TPS) substrate classification task, the metrics are computed on the 512 TPS sequences (**TTT is applied to each protein individually**) based on the cross-validation schema of the TPS dataset (Samusevich et al., 2024). 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			
	mAP \uparrow	AUROC \uparrow	
TerpeneMiner (Samusevich et al., 2024)	0.805	0.948	
TerpeneMiner + TTT (Ours)	0.811 \pm 0.0011	0.950 \pm 0.0002	

Subcellular localization prediction			
	Accuracy \uparrow	MCC \uparrow	F1-score \uparrow
Light attention (Stärk et al., 2021)	0.627	0.549	0.618
Light attention + TTT (Ours)	0.634 \pm 0.004	0.557 \pm 0.005	0.627 \pm 0.004

Results. TTT improves the performance of the base models on both protein function prediction tasks and across all considered metrics (Table 3). Figure 5 provides a qualitative result, where TTT fine-tuning iteratively refines the prediction of TerpeneMiner toward a correct TPS substrate class.

5 DISCUSSION

In this work, we have developed test-time training (TTT) for proteins, enabling per-protein adaptation of machine learning models for enhanced generalization. TTT improves performance across models, their scales, and benchmarks, while primarily enhancing performance on challenging targets. Our results open up the field of self-supervised adaptation for proteins and provide a proof-of-concept for other biology-related domains. While our method demonstrated strong potential, addressing several limitations and researching underexplored directions remain important tasks for future research. Specifically, the success and failure modes of TTT remain unclear, and applying TTT to new tasks requires tuning task-specific hyperparameters. However, our results show that reliable confidence estimates, such as pLDDT, make TTT relatively robust to hyperparameter choices (Figure 13 in Appendix B). Therefore, our future work aims to develop task-agnostic confidence estimates based on protein model representations (Zhang et al., 2024; Rives et al., 2021). Additionally, our findings encourage exploring broader adaptation frameworks for proteins, such as domain adaptation, which leverages both training and test data to address new domains (Ganin & Lempitsky, 2015), and adaptive risk minimization, which employs meta-learning for domain shift adaptation (Zhang et al., 2021).

540 REPRODUCIBILITY STATEMENT

541
542 Our efforts are focused on ensuring that this research is easily reproducible. The proposed test-time
543 training (TTT) method will be released as a Python package, providing easy-to-use wrappers for the
544 models adapted in this paper. Detailed explanations of the application of TTT to individual models
545 and the construction of datasets are included in the appendix. Where applicable, we will also release
546 the source code for dataset generation.

547 ACKNOWLEDGMENTS

548
549 This work was supported by the Ministry of Education, Youth and Sports of the Czech Republic
550 through projects e-INFRA CZ [ID:90254], ELIXIR [LM2023055], CETOCOEN Excellence
551 CZ.02.1.01/0.0/0.0/17_043/0009632, ESFRI RECETOX RI LM2023069. This work was also supported
552 by the European Union (ERC project FRONTIER no. 101097822) and the CETOCOEN
553 EXCELLENCE Teaming project supported from the European Union’s Horizon 2020 research and
554 innovation programme under grant agreement No 857560. This work was also supported by the
555 Czech Science Foundation (GA CR) grant 21-11563M and by the European Union’s Horizon 2020
556 research and innovation programme under Marie Skłodowska-Curie grant agreement No. 891397.
557 Views and opinions expressed are however those of the author(s) only and do not necessarily reflect
558 those of the European Union or the European Research Council. Neither the European Union nor the
559 granting authority can be held responsible for them.

560 REFERENCES

- 561
562 Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf
563 Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure
564 prediction of biomolecular interactions with alphafold 3. *Nature*, pp. 1–3, 2024.
- 565
566 Mehmet Akdel, Douglas EV Pires, Eduard Porta Pardo, Jürgen Jänes, Arthur O Zalevsky, Bálint
567 Mészáros, Patrick Bryant, Lydia L Good, Roman A Laskowski, Gabriele Pozzati, et al. A structural
568 biology community assessment of alphafold2 applications. *Nature Structural & Molecular Biology*,
569 29(11):1056–1067, 2022.
- 570
571 Ethan C Alley, Grigory Khimulya, Surojit Biswas, Mohammed AlQuraishi, and George M Church.
572 Unified rational protein engineering with sequence-based deep representation learning. *Nature*
573 *methods*, 16(12):1315–1322, 2019.
- 574
575 José Juan Almagro Armenteros, Casper Kaae Sønderby, Søren Kaae Sønderby, Henrik Nielsen,
576 and Ole Winther. Deeploc: prediction of protein subcellular localization using deep learning.
577 *Bioinformatics*, 33(21):3387–3395, 2017.
- 578
579 Frances M. Ashcroft, Matthew Lloyd, and Elizabeth A. Haythorne. Glucokinase activity in diabetes:
580 too much of a good thing? *Trends in Endocrinology & Metabolism*, 34(2):119–130, Feb 2023.
581 ISSN 1043-2760. doi: 10.1016/j.tem.2022.12.007. URL <https://doi.org/10.1016/j.tem.2022.12.007>.
- 582
583 Pratyay Banerjee, Tejas Gokhale, and Chitta Baral. Self-supervised test-time learning for reading
584 comprehension. *arXiv preprint arXiv:2103.11263*, 2021.
- 585
586 Inigo Barrio-Hernandez, Jingsi Yeo, Jürgen Jänes, Milot Mirdita, Cameron L. M. Gilchrist, Tanita
587 Wein, Mihaly Varadi, Sameer Velankar, Pedro Beltrao, and Martin Steinegger. Clustering predicted
588 structures at the scale of the known protein universe. *Nature*, 622(7983):637–645, Oct 2023.
589 ISSN 1476-4687. doi: 10.1038/s41586-023-06510-w. URL <https://doi.org/10.1038/s41586-023-06510-w>.
- 590
591 Eyal Ben-David, Nadav Oved, and Roi Reichart. Pada: Example-based prompt learning for on-the-fly
592 adaptation to unseen domains. *Transactions of the Association for Computational Linguistics*, 10:
593 414–433, 2022.
- Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.

- 594 Anton Bushuiev, Roman Bushuiev, Anatolii Filkin, Petr Kouba, Marketa Gabrielova, Michal Gabriel,
595 Jiri Sedlar, Tomas Pluskal, Jiri Damborsky, Stanislav Mazurenko, et al. Learning to design
596 protein-protein interactions with enhanced generalization. *arXiv preprint arXiv:2310.18515*, 2023.
597
- 598 Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony
599 Robinson. One billion word benchmark for measuring progress in statistical language modeling.
600 *arXiv preprint arXiv:1312.3005*, 2013.
- 601 Tianlong Chen and Chengyue Gong. Hotprotein: A novel framework for protein thermostability
602 prediction and editing. *NeurIPS 2022*, 2022.
603
- 604 Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Taylor Applebaum, Alexander
605 Pritzel, Lai Hong Wong, Michal Zielinski, Tobias Sargeant, et al. Accurate proteome-wide
606 missense variant effect prediction with alphamissense. *Science*, 381(6664):eadg7492, 2023.
- 607 Zhixiang Chi, Li Gu, Tao Zhong, Huan Liu, Yuanhao Yu, Konstantinos N Plataniotis, and Yang
608 Wang. Adapting to distribution shift by visual domain prompt generation. *arXiv preprint*
609 *arXiv:2405.02797*, 2024.
610
- 611 David W. Christianson. Structural and chemical biology of terpenoid cyclases. *Chemical Reviews*,
612 117(17):11570–11648, Sep 2017. ISSN 0009-2665. doi: 10.1021/acs.chemrev.7b00287. URL
613 <https://doi.org/10.1021/acs.chemrev.7b00287>.
- 614 The UniProt Consortium. Uniprot: the universal protein knowledgebase in 2023. *Nucleic acids*
615 *research*, 51(D1):D523–D531, 2023.
616
- 617 Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv*
618 *preprint arXiv:1810.04805*, 2018.
619
- 620 Daniel J Diaz, Chengyue Gong, Jeffrey Ouyang-Zhang, James M Loy, Jordan Wells, David Yang,
621 Andrew D Ellington, Alex Dimakis, and Adam R Klivans. Stability oracle: a structure-based
622 graph-transformer for identifying stabilizing mutations. *BioRxiv*, pp. 2023–05, 2023.
- 623 Henry Dieckhaus, Michael Brocidiaco, Nicholas Z Randolph, and Brian Kuhlman. Transfer learn-
624 ing to leverage larger datasets for improved prediction of protein stability changes. *Proceedings of*
625 *the National Academy of Sciences*, 121(6):e2314853121, 2024.
626
- 627 Janani Durairaj, Andrew M Waterhouse, Toomas Mets, Tetiana Brodiazhenko, Minhal Abdullah,
628 Gabriel Studer, Gerardo Tauriello, Mehmet Akdel, Antonina Andreeva, Alex Bateman, et al.
629 Uncovering new families and folds in the natural protein universe. *Nature*, 622(7983):646–653,
630 2023.
- 631 Oliver Dutton, Sandro Bottaro, Istvan Redl, Michele Invernizzi, Albert Chung, Carlo Fiscaro, Falk
632 Hoffmann, Stefano Ruschetta, Fabio Aioldi, Louie Henderson, et al. Improving inverse folding
633 models at protein stability prediction without additional training or data. *bioRxiv*, pp. 2024–06,
634 2024.
- 635 Ahmed Elnaggar, Hazem Essam, Wafaa Salah-Eldin, Walid Moustafa, Mohamed Elkerdawy, Char-
636 lotte Rochereau, and Burkhard Rost. Ankh: Optimized protein language model unlocks general-
637 purpose modelling. *arXiv preprint arXiv:2301.06568*, 2023.
638
- 639 Daniel Esposito, Jochen Weile, Jay Shendure, Lea M Starita, Anthony T Papenfuss, Frederick P Roth,
640 Douglas M Fowler, and Alan F Rubin. Mavedb: an open-source platform to distribute and interpret
641 data from multiplexed assays of variant effect. *Genome biology*, 20:1–11, 2019.
- 642 Xiaomin Fang, Fan Wang, Lihang Liu, Jingzhou He, Dayong Lin, Yingfei Xiang, Kunrui Zhu,
643 Xiaonan Zhang, Hua Wu, Hui Li, et al. A method for multiple-sequence-alignment-free protein
644 structure prediction using a protein language model. *Nature Machine Intelligence*, 5(10):1087–
645 1096, 2023.
646
- 647 Tao Feng, Ziqi Gao, Jiaxuan You, Chenyi Zi, Yan Zhou, Chen Zhang, and Jia Li. Deep reinforcement
learning for modelling protein complexes. *arXiv preprint arXiv:2405.02299*, 2024.

- 648 Noelia Ferruz, Steffen Schmidt, and Birte Höcker. Protgpt2 is a deep unsupervised language model
649 for protein design. *Nature communications*, 13(1):4348, 2022.
- 650
- 651 Jonathan Frazer, Pascal Notin, Mafalda Dias, Aidan Gomez, Joseph K Min, Kelly Brock, Yarin Gal,
652 and Debora S Marks. Disease variant prediction with deep generative models of evolutionary data.
653 *Nature*, 599(7883):91–95, 2021.
- 654 Yossi Gandelsman, Yu Sun, Xinlei Chen, and Alexei Efros. Test-time training with masked autoen-
655 coders. *Advances in Neural Information Processing Systems*, 35:29374–29385, 2022.
- 656
- 657 Yaroslav Ganin and Victor S. Lempitsky. Unsupervised domain adaptation by backpropagation. In
658 Francis R. Bach and David M. Blei (eds.), *Proceedings of the 32nd International Conference on*
659 *Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop*
660 *and Conference Proceedings*, pp. 1180–1189. JMLR.org, 2015. URL [http://proceedings.](http://proceedings.mlr.press/v37/ganin15.html)
661 [mlr.press/v37/ganin15.html](http://proceedings.mlr.press/v37/ganin15.html).
- 662 Jan Gorodkin. Comparing two k-category assignments by a k-category correlation coefficient.
663 *Computational biology and chemistry*, 28(5-6):367–374, 2004.
- 664 Xin Gu, Patrick Jouandin, Pranav V. Lalgudi, Rich Binari, Max L. Valenstein, Michael A. Reid,
665 Annamarie E. Allen, Nolan Kamitaki, Jason W. Locasale, Norbert Perrimon, and David M.
666 Sabatini. Sestrin mediates detection of and adaptation to low-leucine diets in drosophila. *Nature*,
667 608(7921):209–216, Aug 2022. ISSN 1476-4687. doi: 10.1038/s41586-022-04960-2. URL
668 <https://doi.org/10.1038/s41586-022-04960-2>.
- 669
- 670 Muhammet F. Gulen, Natasha Samson, Alexander Keller, Marius Schwabenland, Chong Liu, Selene
671 Glück, Vivek V. Thacker, Lucie Favre, Bastien Mangeat, Lona J. Kroese, Paul Krimpenfort, Marco
672 Prinz, and Andrea Ablasser. cgas–sting drives ageing-related inflammation and neurodegeneration.
673 *Nature*, 620(7973):374–380, Aug 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06373-1.
674 URL <https://doi.org/10.1038/s41586-023-06373-1>.
- 675 Kathryn H. Gunn and Saskia B. Neher. Structure of dimeric lipoprotein lipase reveals a pore adjacent
676 to the active site. *Nature Communications*, 14(1):2569, May 2023. ISSN 2041-1723. doi: 10.1038/
677 [s41467-023-38243-9](https://doi.org/10.1038/s41467-023-38243-9). URL <https://doi.org/10.1038/s41467-023-38243-9>.
- 678 Moritz Hardt and Yu Sun. Test-time training on nearest neighbors for large language models. *arXiv*
679 *preprint arXiv:2305.18466*, 2023.
- 680
- 681 Tomas Hayes, Roshan Rao, Halil Akin, Nicholas J Sofroniew, Deniz Oktay, Zeming Lin, Robert
682 Verkuil, Vincent Q Tran, Jonathan Deaton, Marius Wiggert, et al. Simulating 500 million years of
683 evolution with a language model. *bioRxiv*, pp. 2024–07, 2024.
- 684 Michael Heinzinger, Konstantin Weissenow, Joaquin Gomez Sanchez, Adrian Henkel, Milot Mirdita,
685 Martin Steinegger, and Burkhard Rost. Bilingual language model for protein sequence and structure.
686 *bioRxiv*, pp. 2023–07, 2023.
- 687
- 688 Lucas Torroba Hennigen and Yoon Kim. Deriving language models from masked language models.
689 *arXiv preprint arXiv:2305.15501*, 2023.
- 690
- 691 Thomas A Hopf, John B Ingraham, Frank J Poelwijk, Charlotta PI Schärfe, Michael Springer, Chris
692 Sander, and Debora S Marks. Mutation effects predicted from sequence co-variation. *Nature*
693 *biotechnology*, 35(2):128–135, 2017.
- 694
- 695 Gerta Hoxhaj and Brendan D. Manning. The pi3k–akt network at the interface of oncogenic signalling
696 and cancer metabolism. *Nature Reviews Cancer*, 20(2):74–88, Feb 2020. ISSN 1474-1768. doi: 10.
697 [1038/s41568-019-0216-7](https://doi.org/10.1038/s41568-019-0216-7). URL <https://doi.org/10.1038/s41568-019-0216-7>.
- 698
- 699 Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and Alexander
700 Rives. Learning inverse folding from millions of predicted structures. In *International conference*
701 *on machine learning*, pp. 8946–8970. PMLR, 2022.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
arXiv:2106.09685, 2021.

- 702 Bowen Jing, Bonnie Berger, and Tommi Jaakkola. Alphafold meets flow matching for generating
703 protein ensembles. *arXiv preprint arXiv:2402.04845*, 2024.
704
- 705 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,
706 Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate
707 protein structure prediction with alphafold. *nature*, 596(7873):583–589, 2021.
- 708 Pranav Kantroo, Gunter Wagner, and Benjamin Machta. Pseudo-perplexity in one fell swoop for
709 protein fitness estimation. *bioRxiv*, pp. 2024–07, 2024.
- 710 Neerav Karani, Ertunc Erdil, Krishna Chaitanya, and Ender Konukoglu. Test-time adaptable neural
711 networks for robust medical image segmentation. *Medical Image Analysis*, 68:101907, 2021.
712
- 713 Zuzana Keckesova, Joana Liu Donaher, Jasmine De Cock, Elizaveta Freinkman, Susanne Lingrell,
714 Daniel A. Bachovchin, Brian Bierie, Verena Tischler, Aurelia Noske, Marian C. Okondo, Ferenc
715 Reinhardt, Prathapan Thiru, Todd R. Golub, Jean E. Vance, and Robert A. Weinberg. Lactb is a
716 tumour suppressor that modulates lipid metabolism and cell state. *Nature*, 543(7647):681–686,
717 Mar 2017. ISSN 1476-4687. doi: 10.1038/nature21408. URL [https://doi.org/10.1038/
718 nature21408](https://doi.org/10.1038/nature21408).
- 719 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua
720 Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR
721 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL [http:
722 //arxiv.org/abs/1412.6980](http://arxiv.org/abs/1412.6980).
- 723 Andrew Kirjner, Jason Yim, Raman Samusevich, Shahar Bracha, Tommi S Jaakkola, Regina Barzilay,
724 and Ila R Fiete. Improving protein optimization with smoothed fitness landscapes. In *The Twelfth
725 International Conference on Learning Representations*, 2023.
726
- 727 Petr Kouba, Pavel Kohout, Faraneh Haddadi, Anton Bushuiev, Raman Samusevich, Jiri Sedlar, Jiri
728 Damborsky, Tomas Pluskal, Josef Sivic, and Stanislav Mazurenko. Machine learning-guided
729 protein engineering. *ACS catalysis*, 13(21):13863–13895, 2023.
- 730 Andriy Kryshtafovych, Torsten Schwede, Maya Topf, Krzysztof Fidelis, and John Moulton. Critical
731 assessment of methods of protein structure prediction (casp)—round xv. *Proteins: Structure,
732 Function, and Bioinformatics*, 91(12):1539–1549, 2023. doi: <https://doi.org/10.1002/prot.26617>.
733 URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/prot.26617>.
- 734 Elodie Laine, Yasaman Karami, and Alessandra Carbone. Gemme: a simple and fast global epistatic
735 model predicting mutational effects. *Molecular biology and evolution*, 36(11):2604–2619, 2019.
736
- 737 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,
738 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom
739 Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level pro-
740 tein structure with a language model. *Science*, 379(6637):1123–1130, 2023. doi: 10.1126/
741 science.ade2574. URL [https://www.science.org/doi/abs/10.1126/science.
742 ade2574](https://www.science.org/doi/abs/10.1126/science.ade2574).
- 743 Yuejiang Liu, Parth Kothari, Bastien van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and Alexandre
744 Alahi. TTT++: when does self-supervised test-time training fail or thrive? In Marc’Aurelio
745 Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan
746 (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neu-
747 ral Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp.
748 21808–21820, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/
749 b618c3210e934362ac261db280128c22-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/b618c3210e934362ac261db280128c22-Abstract.html).
- 750 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International
751 Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*.
752 OpenReview.net, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- 753 Ali Madani, Ben Krause, Eric R Greene, Subu Subramanian, Benjamin P Mohr, James M Holton,
754 Jose Luis Olmos, Caiming Xiong, Zachary Z Sun, Richard Socher, et al. Large language models
755 generate functional protein sequences across diverse families. *Nature Biotechnology*, 41(8):
1099–1106, 2023.

- 756 Valerio Mariani, Marco Biasini, Alessandro Barbato, and Torsten Schwede. Iddt: a local superposition-
757 free score for comparing protein structures and models using distance difference tests. *Bioinform-*
758 *atics*, 29(21):2722–2728, 2013.
- 759 Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, and Alex Rives. Language models
760 enable zero-shot prediction of the effects of mutations on protein function. *Advances in neural*
761 *information processing systems*, 34:29287–29303, 2021.
- 762 Peter Mikhael, Itamar Chinn, and Regina Barzilay. Clipzyme: Reaction-conditioned virtual screening
763 of enzymes. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Aus-*
764 *tria, July 21-27, 2024*. OpenReview.net, 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=0mYAK6Yhhm)
765 [id=0mYAK6Yhhm](https://openreview.net/forum?id=0mYAK6Yhhm).
- 766 Amir Motmaen, Justas Dauparas, Minkyung Baek, Mohamad H Abedi, David Baker, and Philip
767 Bradley. Peptide-binding specificity prediction using fine-tuned protein structure prediction
768 networks. *Proceedings of the National Academy of Sciences*, 120(9):e2216697120, 2023.
- 769 Pascal Notin, Mafalda Dias, Jonathan Frazer, Javier Marchena-Hurtado, Aidan N Gomez, Debora
770 Marks, and Yarin Gal. Tranception: protein fitness prediction with autoregressive transformers
771 and inference-time retrieval. In *International Conference on Machine Learning*, pp. 16990–17017.
772 PMLR, 2022a.
- 773 Pascal Notin, Lood Van Niekerk, Aaron W Kollasch, Daniel Ritter, Yarin Gal, and Debora S Marks.
774 Trancepteve: Combining family-specific and family-agnostic models of protein sequences for
775 improved fitness prediction. *bioRxiv*, pp. 2022–12, 2022b.
- 776 Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood Van Niekerk, Steffanie Paul, Han Spinner, Nathan
777 Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, et al. Proteingym: Large-scale benchmarks
778 for protein fitness prediction and design. *Advances in Neural Information Processing Systems*, 36,
779 2024.
- 780 Dong oh Seo, David O’Donnell, Nimansha Jain, Jason D. Ulrich, Jasmin Herz, Yuhao Li, Mackenzie
781 Lemieux, Jiye Cheng, Hao Hu, Javier R. Serrano, Xin Bao, Emily Franke, Maria Karlsson,
782 Martin Meier, Su Deng, Chandani Desai, Hemraj Dodiya, Janaki Lelwala-Guruge, Scott A.
783 Handley, Jonathan Kipnis, Sangram S. Sisodia, Jeffrey I. Gordon, and David M. Holtzman.
784 Apoe isoform- and microbiota-dependent progression of neurodegeneration in a mouse model of
785 tauopathy. *Science*, 379(6628):eadd1236, 2023. doi: 10.1126/science.add1236. URL <https://www.science.org/doi/abs/10.1126/science.add1236>.
- 786 Andrei Papkou, Lucia Garcia-Pastor, José Antonio Escudero, and Andreas Wagner. A rugged yet eas-
787 ily navigable fitness landscape. *Science*, 382(6673):eadh3860, 2023. doi: 10.1126/science.adh3860.
788 URL <https://www.science.org/doi/abs/10.1126/science.adh3860>.
- 789 Predrag Radivojac and et al. A large-scale evaluation of computational protein function prediction.
790 *Nature Methods*, 10(3):221–227, Mar 2013. ISSN 1548-7105. doi: 10.1038/nmeth.2340. URL
791 <https://doi.org/10.1038/nmeth.2340>.
- 792 Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter Abbeel,
793 and Yun Song. Evaluating protein transfer learning with tape. *Advances in neural information*
794 *processing systems*, 32, 2019.
- 795 Roshan Rao, Joshua Meier, Tom Sercu, Sergey Ovchinnikov, and Alexander Rives. Transformer
796 protein language models are unsupervised structure learners. *Biorxiv*, pp. 2020–12, 2020.
- 797 Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu,
798 and Alexander Rives. Msa transformer. In *International Conference on Machine Learning*, pp.
799 8844–8856. PMLR, 2021.
- 800 Ali Razavi, Aäron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity im-
801 ages with VQ-VAE-2. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Flo-
802 rence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural In-*
803 *formation Processing Systems 32: Annual Conference on Neural Information Process-*
804 *ing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp.
805 806 807 808 809

- 810 14837–14847, 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/](https://proceedings.neurips.cc/paper/2019/hash/5f8e2fa1718d1bbcadf1cd9c7a54fb8c-Abstract.html)
811 [5f8e2fa1718d1bbcadf1cd9c7a54fb8c-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/5f8e2fa1718d1bbcadf1cd9c7a54fb8c-Abstract.html).
812
- 813 Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo,
814 Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from
815 scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National*
816 *Academy of Sciences*, 118(15):e2016239118, 2021.
- 817 Xavier Robin, Juergen Haas, Rafal Gumienny, Anna Smolinski, Gerardo Tauriello, and Torsten
818 Schwede. Continuous automated model evaluation (cameo)—perspectives on the future of fully
819 automated evaluation of structure prediction methods. *Proteins: Structure, Function, and Bioinfor-*
820 *matics*, 89(12):1977–1986, 2021.
821
- 822 Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint*
823 *arXiv:1609.04747*, 2016.
- 824 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.
825 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceed-*
826 *ings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22500–22510,
827 2023.
828
- 829 Julian Salazar, Davis Liang, Toan Q Nguyen, and Katrin Kirchhoff. Masked language model scoring.
830 *arXiv preprint arXiv:1910.14659*, 2019.
- 831 Raman Samusevich, Téo Hebra, Roman Bushuiev, Anton Bushuiev, Tereza Čalounová, Helena
832 Smrčková, Ratthachat Chatpatanasiri, Jonáš Kulháněk, Milana Perković, Martin Engst, Adéla
833 Tajovská, Josef Sivic, and Tomáš Pluskal. Highly accurate discovery of terpene synthases powered
834 by machine learning reveals functional terpene cyclization in archaea. *bioRxiv*, 2024. doi: 10.
835 1101/2024.01.29.577750. URL [https://www.biorxiv.org/content/early/2024/](https://www.biorxiv.org/content/early/2024/04/25/2024.01.29.577750)
836 [04/25/2024.01.29.577750](https://www.biorxiv.org/content/early/2024/04/25/2024.01.29.577750).
837
- 838 Nicolae Sapoval, Amirali Aghazadeh, Michael G. Nute, Dinler A. Antunes, Advait Balaji, Richard
839 Baraniuk, C. J. Barberan, Ruth Dannenfeler, Chen Dun, Mohammadamin Edrisi, R. A. Leo
840 Elworth, Bryce Kille, Anastasios Kyrillidis, Luay Nakhleh, Cameron R. Wolfe, Zhi Yan, Vicky
841 Yao, and Todd J. Treangen. Current progress and open challenges for applying deep learning across
842 the biosciences. *Nature Communications*, 13(1):1728, Apr 2022. ISSN 2041-1723. doi: 10.1038/
843 [s41467-022-29268-7](https://doi.org/10.1038/s41467-022-29268-7). URL <https://doi.org/10.1038/s41467-022-29268-7>.
- 844 Emre Sevgen, Joshua Moller, Adrian Lange, John Parker, Sean Quigley, Jeff Mayer, Poonam
845 Srivastava, Sitaram Gayatri, David Hosfield, Maria Korshunova, et al. Prot-vae: protein transformer
846 variational autoencoder for functional protein design. *bioRxiv*, pp. 2023–01, 2023.
847
- 848 Raghav Shroff, Austin W Cole, Daniel J Diaz, Barrett R Morrow, Isaac Donnell, Ankur Annappareddy,
849 Jimmy Gollihar, Andrew D Ellington, and Ross Thyer. Discovery of novel gain-of-function
850 mutations guided by structure-based deep learning. *ACS synthetic biology*, 9(11):2927–2935, 2020.
- 851 Yidong Song, Qianmu Yuan, Sheng Chen, Yuansong Zeng, Huiying Zhao, and Yuedong Yang.
852 Accurately predicting enzyme functions through geometric graph learning on esmfold-predicted
853 structures. *Nature Communications*, 15(1):8180, 2024.
854
- 855 Hannes Stärk, Christian Dallago, Michael Heinzinger, and Burkhard Rost. Light attention predicts
856 protein location from the language of life. *Bioinformatics Advances*, 1(1):vbab035, 11 2021. ISSN
857 2635-0041. doi: 10.1093/bioadv/vbab035. URL [https://doi.org/10.1093/bioadv/](https://doi.org/10.1093/bioadv/vbab035)
858 [vbab035](https://doi.org/10.1093/bioadv/vbab035).
- 859 Jin Su, Chenchen Han, Yuyang Zhou, Junjie Shan, Xibin Zhou, and Fajie Yuan. Saprot: Protein
860 language modeling with structure-aware vocabulary. *bioRxiv*, pp. 2023–10, 2023.
861
- 862 Sriram Subramaniam and Gerard J. Kleywegt. A paradigm shift in structural biology. *Nature*
863 *Methods*, 19(1):20–23, Jan 2022. ISSN 1548-7105. doi: 10.1038/s41592-021-01361-7. URL
<https://doi.org/10.1038/s41592-021-01361-7>.

- 864 Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training
865 with self-supervision for generalization under distribution shifts. In *International conference on*
866 *machine learning*, pp. 9229–9248. PMLR, 2020.
- 867
868 Nataša Tagasovska, Ji Won Park, Matthieu Kirchmeyer, Nathan C Frey, Andrew Martin Watkins,
869 Aya Abdelsalam Ismail, Arian Rokkum Jamasb, Edith Lee, Tyler Bryson, Stephen Ra, et al.
870 Antibody domainbed: Out-of-distribution generalization in therapeutic protein design. *arXiv*
871 *preprint arXiv:2407.21028*, 2024.
- 872 Kotaro Tsuboyama, Justas Dauparas, Jonathan Chen, Elodie Laine, Yasser Mohseni Behbahani,
873 Jonathan J Weinstein, Niall M Mangan, Sergey Ovchinnikov, and Gabriel J Rocklin. Mega-scale
874 experimental analysis of protein folding stability in biology and design. *Nature*, 620(7973):
875 434–444, 2023.
- 876 Mike Tyers and Matthias Mann. From genomics to proteomics. *Nature*, 422(6928):193–197, Mar
877 2003. ISSN 1476-4687. doi: 10.1038/nature01510. URL [https://doi.org/10.1038/
878 nature01510](https://doi.org/10.1038/nature01510).
- 879
880 Michel van Kempen, Stephanie S Kim, Charlotte Tumescheit, Milot Mirdita, Cameron LM Gilchrist,
881 Johannes Söding, and Martin Steinegger. Foldseek: fast and accurate protein structure search.
882 *Biorxiv*, pp. 2022–02, 2022.
- 883
884 Mihaly Varadi, Stephen Anyango, Mandar Deshpande, Sreenath Nair, Cindy Natassia, Galabina
885 Yordanova, David Yuan, Oana Stroe, Gemma Wood, Agata Laydon, et al. Alphafold protein
886 structure database: massively expanding the structural coverage of protein-sequence space with
887 high-accuracy models. *Nucleic acids research*, 50(D1):D439–D444, 2022.
- 888 A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- 889
890 Renhao Wang, Yu Sun, Yossi Gandelsman, Xinlei Chen, Alexei A Efros, and Xiaolong Wang.
891 Test-time training on video streams. *arXiv preprint arXiv:2307.05014*, 2023.
- 892 Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach,
893 Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein
894 structure and function with rfdiffusion. *Nature*, 620(7976):1089–1100, 2023.
- 895
896 Zehao Xiao, Xiantong Zhen, Ling Shao, and Cees GM Snoek. Learning to generalize across domains
897 on single test samples. *arXiv preprint arXiv:2202.08045*, 2022.
- 898 Jason Yang, Francesca-Zhoufan Li, and Frances H. Arnold. Opportunities and challenges for machine
899 learning-assisted enzyme engineering. *ACS Central Science*, 10(2):226–241, Feb 2024. ISSN 2374-
900 7943. doi: 10.1021/acscentsci.3c01275. URL [https://doi.org/10.1021/acscentsci.
901 3c01275](https://doi.org/10.1021/acscentsci.3c01275).
- 902
903 Zhenyu Yang, Xiaoxi Zeng, Yi Zhao, and Runsheng Chen. Alphafold2 and its applications in the
904 fields of biology and medicine. *Signal Transduction and Targeted Therapy*, 8(1):115, 2023.
- 905
906 Tianhao Yu, Haiyang Cui, Jianan Canal Li, Yunan Luo, Guangde Jiang, and Huimin Zhao. En-
907 zyme function prediction using contrastive learning. *Science*, 379(6639):1358–1363, 2023.
908 doi: 10.1126/science.adf2465. URL [https://www.science.org/doi/abs/10.1126/
909 science.adf2465](https://www.science.org/doi/abs/10.1126/science.adf2465).
- 910 Marvin Zhang, Henrik Marklund, Nikita Dhawan, Abhishek Gupta, Sergey Levine, and Chelsea
911 Finn. Adaptive risk minimization: Learning to adapt to domain shift. In Marc’Aurelio Ran-
912 zato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan
913 (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neu-
914 ral Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp.
915 23664–23678, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/
916 c705112d1ec18b97acac7e2d63973424-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/c705112d1ec18b97acac7e2d63973424-Abstract.html).
- 917
918 Yang Zhang and Jeffrey Skolnick. Scoring function for automated assessment of protein structure
919 template quality. *Proteins: Structure, Function, and Bioinformatics*, 57(4):702–710, 2004.

918 Zhidian Zhang, Hannah K. Wayment-Steele, Garyk Brixii, Haobo Wang, Dorothee Kern, and Sergey
919 Ovchinnikov. Protein language models learn evolutionary statistics of interacting sequence
920 motifs. *Proceedings of the National Academy of Sciences*, 121(45):e2406285121, 2024. doi:
921 10.1073/pnas.2406285121. URL [https://www.pnas.org/doi/abs/10.1073/pnas.](https://www.pnas.org/doi/abs/10.1073/pnas.2406285121)
922 2406285121.

923 Hao Zhao, Yuejiang Liu, Alexandre Alahi, and Tao Lin. On pitfalls of test-time adaptation. In Andreas
924 Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett
925 (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu,*
926 *Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 42058–42080.
927 PMLR, 2023. URL <https://proceedings.mlr.press/v202/zhao23d.html>.

928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

APPENDIX

In Appendix A, we provide further details on the experimental setup, including comprehensive descriptions of the models, datasets, and metrics used. Next, in Appendix B, we present additional results and their analysis. We discuss the distribution of TTT effects and demonstrate that TTT primarily improves performance on challenging targets. We also explore the impact of hyperparameters by showing the performance on validation sets.

A EXPERIMENTAL DETAILS

In this section, we describe the experimental details for the three downstream tasks considered in this work: protein fitness prediction (Appendix A.1), protein structure prediction (Appendix A.2), and protein function prediction (Appendix A.3). Each subsection describes the application of test-time training (TTT) to the respective models, along with details on the datasets, metrics, and models. Table 4 additionally summarizes the hyperparameters used for the application of TTT to individual models.

A.1 PROTEIN FITNESS PREDICTION

A.1.1 DATASETS

ProteinGym. ProteinGym³ is the standard benchmark for protein fitness prediction (Notin et al., 2024). The latest, second version of the dataset includes 217 deep mutation scanning experiments (DMSs) across different proteins. **We focus on the well-established zero-shot variant of the benchmark and do not experiment with the supervised variant, as it has not yet been fully incorporated into the official codebase at the time of this study.** In total, the dataset contains 2.5 mutants with annotated 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 MaveDB⁴ (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.

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 6A). To ensure the quality of the constructed MaveDB dataset, we validated that representative baselines from ProteinGym generalize to the new assays, following a similar distribution of predictions (Figure 6B,C). Finally, for efficiently tuning hyper-parameters for fitness prediction models we sampled 50 random proteins (Figure 6D), corresponding to 83 assays and collectively 134 thousand variants.

A.1.2 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.

³<https://github.com/OATML-Markslab/ProteinGym>

⁴<https://www.mavedb.org>

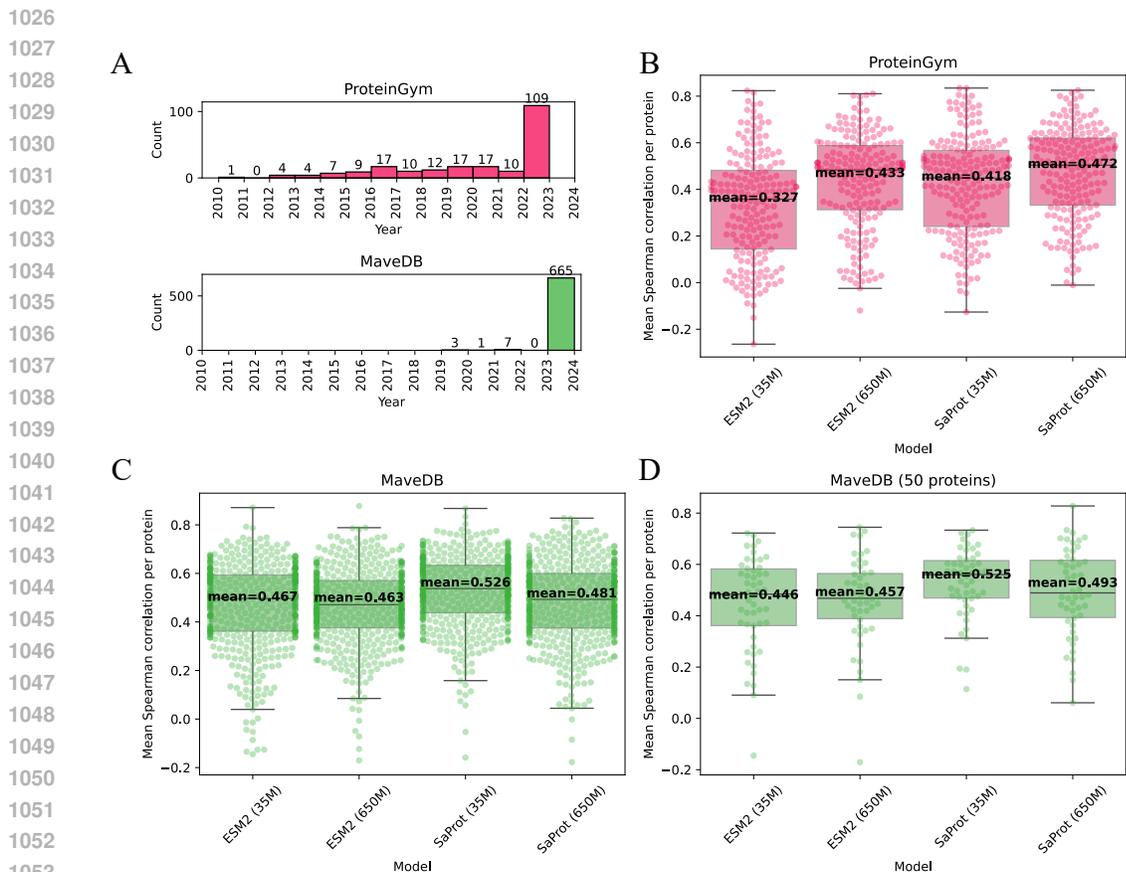


Figure 6: Comparison of the standard ProteinGym dataset with the MaveDB dataset constructed in this work. A) MaveDB, mined from Esposito et al. (2019), includes novel assays even after filtering to ensure distinct proteins from the comprehensive ProteinGym dataset. This is largely because most MaveDB 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.

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 introduced mutations: protein 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 1, Table 5).

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 6). 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.

A.1.3 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 protein language models 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:

$$\sum_{i \in T} \log p(x_i^m | x_{\setminus i}) - \log p(x_i | x_{\setminus i}), \quad (3)$$

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 + TTT. ESM2 can be straightforwardly enhanced with test-time training. 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 (3) serves as the task-specific head h , which in this case involves the pre-training head g that predicts log probabilities. Overall, we apply TTT to the pre-trained ESM2 model and, after a pre-defined number of self-supervised fine-tuning steps, score mutations using Equation (3). During TTT we fine-tune all parameters in $g \circ f$ end-to-end except for token and position embeddings.

SaProt. We also experiment with the state-of-the-art fitness prediction 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 (3) 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 (3) are replaced with $x_{\setminus T}$ instead of $x_{\setminus i}$, not assuming the independence of substitutions. During TTT, 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 ProteinGym³.

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

A.2 PROTEIN STRUCTURE PREDICTION

A.2.1 DATASETS

CAMEO dataset. To evaluate the capabilities of TTT on protein folding, we employ the CAMEO validation and test sets as described in Lin et al. (2023). Specifically, the validation set was obtained

1134 by querying the CAMEO (Continuous Automated Model Evaluation) web server⁵ (Robin et al., 2021)
1135 for entries between August 2021 and January 2022, while the CAMEO test set consists of entries
1136 from April 1, 2022, to June 25, 2022. Most of the entries in the CAMEO sets are predicted with high
1137 accuracy and confidence (Lin et al., 2023). Therefore, we subselected the challenging validation and
1138 test sets where TTT is relevant.

1139 Specifically, we applied two criteria: (1) preserving entries with ESMFold pLDDT scores below 70
1140 to filter out high-confidence predictions (Jumper et al., 2021), and (2) selecting entries with ESM2
1141 perplexity scores greater than or equal to 6, ensuring that the predictions are challenging due to poor
1142 sequence understanding rather than other factors. Additionally, most structures with perplexity scores
1143 below 6 are already associated with high-confidence predictions (Figure S5 in Lin et al. (2023)).
1144 After filtering, the resulting challenging validation and test sets consist of 27 (out of 378) and 18 (out
1145 of 194) targets, respectively. The vast majority of the remaining structures have accurate ESMFold
1146 structure predictions.

1147 A.2.2 METRICS

1149 To assess the quality of the predicted protein structures with respect to the ground truth structures, we
1150 use two standard metrics averaged across the test dataset: TM-score (Zhang & Skolnick, 2004) and
1151 LDDT (Mariani et al., 2013).

1153 **TM-score.** The TM-score (Template Modeling score) is a metric used to assess the quality of the
1154 global 3D alignment between the predicted and target protein structures. It evaluates the structural
1155 similarity by comparing the distance between corresponding residues after superposition. The
1156 TM-score ranges from 0 to 1, where higher values indicate better alignment.

1158 **LDDT.** The Local Distance Difference Test (LDDT) is an alignment-free metric used to assess the
1159 accuracy of predicted protein structures. Unlike global metrics, LDDT focuses on local structural
1160 differences by measuring the deviation in distances between atom pairs in the predicted structure
1161 compared to the target structure. It is particularly useful for evaluating the accuracy of local regions,
1162 such as secondary structure elements. LDDT scores range from 0 to 100, with higher values indicating
1163 better local structural agreement.

1164 A.2.3 MODELS

1166 **ESMFold.** The ESMFold architecture comprises two key components: a protein language model,
1167 ESM2, which, given a protein sequence, generates embeddings for individual amino acids, and a
1168 folding block that, using these embeddings and the sequence, predicts the protein 3D structure along
1169 with per-amino-acid confidence scores, known as pLDDT scores. In our experiments, we use the
1170 `esmfold_v0` model from the publicly available ESMFold checkpoints⁶. Please note that we use
1171 `esmfold_v0` and not `esmfold_v1` to avoid data leakage with respect to the CAMEO test set.

1173 **ESMFold + TTT.** Since ESM2 backbone of ESMFold was pre-trained in a self-supervised masked
1174 modeling regime, the application of TTT to ESMFold is straightforward. We treat ESM2 as the
1175 backbone f , the language modeling head predicting amino acid classes from their embeddings as the
1176 self-supervised head g , and the folding trunk along with the structure modules as the downstream
1177 task head h . After each TTT step, we run $h \circ f$ to compute the pLDDT scores, which allows us to
1178 estimate the optimal number of TTT steps for each protein based on the highest pLDDT score.

1179 Since the backbone f is given by the ESM2 model containing 3 billion parameters, we apply LoRA
1180 (Hu et al., 2021) to all matrices involved in self-attention. This enables fine-tuning ESMFold + TTT
1181 on a single GPU.

1182 **ESMFold + ME.** Since ESMFold is a regression model, it only predicts one solution and does
1183 not have a straightforward mechanism of sampling multiple structure predictions. Nevertheless, the
1184 authors of ESMFold propose a way to sample multiple candidates (Section A.3.2 in Lin et al. (2023)).
1185

1186 ⁵<https://www.cameo3d.org/modeling>

1187 ⁶<https://github.com/facebookresearch/esm/blob/main/esm/esmfold/v1/pretrained.py>

To sample more solutions, the masking prediction (ME) method randomly masks 15% (same ratio as during masked language modeling pre-training) of the amino acids embeddings before passing them to the structure prediction block. Selecting the solution with the highest pLDDT may lead to improved predicted structure. Since sampling multiple solutions with ESMFold + ME and selecting the best one via pLDDT is analogous to ESMFold + TTT, we employ the former as a baseline, running the method for the same number of step.

ESM3. Unlike ESMFold, ESM3 is a fully multiple-track, BERT-like model (Devlin, 2018), pre-trained to unmask both protein sequence and structure tokens simultaneously (along with the function tokens). The structure tokens in ESM3 are generated via a separately pre-trained VQ-VAE (Razavi et al., 2019) operating on the protein geometry. In our experiments, we use the smallest, publicly available version of the ESM3 model (ESM3_sm_open_v0)⁷.

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

Despite the fact that the model contains 1.4 billion parameters, even without using LoRA, ESM3 + TTT can be fine-tuned on a single NVIDIA A100 GPU. Therefore, we do not employ LoRA for fine-tuning ESM3.

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

A.3 PROTEIN FUNCTION PREDICTION

A.3.1 DATASETS

TPS dataset. For the evaluation of terpene substrate classification, we use the largest available dataset of characterized TPS enzymes from Samusevich et al. (2024) and repurpose the original 5-fold cross-validation schema. We focus on the most challenging TPS sequences, defined as those predicted by the TPS detector, proposed by the dataset authors, with confidence scores below 0.8. This filtering results in 104, 98, 113, 100, 97 examples in the individual folds.

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

DeepLoc. For hyperparameter tuning in the subcellular location prediction task, we use the test set from the DeepLoc dataset (Almagro Armenteros et al., 2017). Similar to setHard, DeepLoc assigns labels from one of ten subcellular location classes. The dataset contains 2768 proteins, which we further filter to 2457 sequences that do not exceed a length of 1022 amino acids, ensuring

⁷<https://github.com/evolutionaryscale/esm>

1242 compatibility with the embedding capabilities of ESM-1b. Since setHard was constructed to be
1243 independent of DeepLoc, setHard provides a leakage-free source of data for validation.
1244

1245 A.3.2 METRICS

1246
1247 **mAP, AUROC.** The TPS substrate prediction problem is a 12-class multi-label classification task
1248 over possible TPS substrates. Therefore, we assess the quality of the predictions using standard
1249 multi-label classification metrics such as mean average precision (mAP) and area under the receiver
1250 operating characteristic curve (AUROC) averaged across individual classes. These metrics were
1251 used in the original TerpeneMiner paper (Samusevich et al., 2024). We report the performance by
1252 averaging the metric values concatenated across all validation folds from the 5-fold cross-validation
1253 schema.

1254 **Accuracy, MCC, F1-score.** To evaluate the performance of subcellular location prediction methods,
1255 we use standard classification metrics as employed in Stärk et al. (2021). Accuracy standardly
1256 measures the ratio of correctly classified proteins, while Matthew’s correlation coefficient for multiple
1257 classes (MCC) serves as an alternative to the Pearson correlation coefficient for classification tasks
1258 (Gorodkin, 2004). The F1-score, the harmonic mean of precision and recall, evaluates performance
1259 from a retrieval perspective, balancing the trade-off between false positives and false negatives.
1260

1261 A.3.3 MODELS

1262
1263 **TerpeneMiner.** TerpeneMiner is a state-of-the-art method for the classification of terpene synthase
1264 (TPS) substrates (Samusevich et al., 2024). The model consists of two parallel tracks. Given a
1265 protein sequence, TerpeneMiner first computes its ESM-1v embedding (Meier et al., 2021) and
1266 a vector of similarities to the functional domains of proteins from the training dataset, based on
1267 unsupervised domain segmentation of AlphaFold2-predicted structures (Jumper et al., 2021). The
1268 ESM-1v embedding and the similarity vector are then concatenated and processed by a separately
1269 trained random forest, which predicts TPS substrate class probabilities.

1270 In our experiments, we use the “PLM only” version of the model, which leverages only ESM-1v
1271 embeddings (PLM stands for protein language model). This version exhibits a minor performance de-
1272 crease compared to the full model but exactly follows a Y-shaped architecture, allowing us to validate
1273 the effectiveness of test-time training for predicting TPS substrates. We use the implementation of
1274 TerpeneMiner available at the official GitHub page⁸.

1275 **TerpeneMiner + TTT.** When applying TTT to TerpeneMiner, we treat the frozen ESM-1v model
1276 as a backbone f , its language modeling head as a self-supervised head g , and the random forest
1277 classifying TPS substrates as a downstream supervised head h .
1278

1279 **Light Attention.** We use Light attention (Stärk et al., 2021) as a representative baseline for
1280 subcellular location prediction. Light attention leverages protein embeddings from a language model,
1281 which in our case is ESM-1b (Rives et al., 2021). The model processes per-residue embeddings via a
1282 softmax-weighted aggregation mechanism, referred to as light attention, which operates with linear
1283 complexity relative to sequence length and enables richer aggregation of per-residue information, as
1284 opposed to standard mean pooling. We re-train the model using ESM-1b embeddings on the DeepLoc
1285 dataset (Almagro Armenteros et al., 2017) using the code from the official GitHub page⁹.
1286

1287 **Light attention + TTT.** When applying TTT to Light attention, we treat the frozen ESM-1b as the
1288 backbone f , the language modeling head of ESM-1b as the self-supervised head g , and the Light
1289 attention block as the fine-tuning head h .

1290 B EXTENDED RESULTS

1291
1292 In this section, we provide additional results on test sets (Appendix B.1) and discuss validation
1293 performance (Appendix B.2).
1294

1295 ⁸<https://github.com/pluskal-lab/TerpeneMiner>

⁹<https://github.com/HannesStark/protein-localization>

Table 4: **Hyperparameters used for adapting TTT 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 TTT steps, as adjusting the learning rate and batch size effectively controls the expected performance under the fixed number of steps, as shown in Figure 12. Therefore, we used 30 steps in all our experiments. The only exception was ESM3 + TTT, 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.**

	Learning rate	Batch size	Grad. acc. steps	TTT steps	LoRA rank r	LoRA α
<i>Fitness prediction</i>						
ESM2 (35M) + TTT	4e-5, 4e-4 , 4e-3	4	4, 8, 16 , 32, 64	30	-, 4 , 8 , 32	-, 8 , 16 , 32
ESM2 (650M) + TTT	4e-5 , 4e-4, 4e-3	4	4, 8, 16 , 32	30	-, 4 , 8 , 32	-, 8 , 16 , 32
SaProt (35M) + TTT	4e-5, 4e-4 , 4e-3	4	4, 8 , 16, 32	30	-	-
SaProt (650M) + TTT	4e-5 , 4e-4, 4e-3	2 , 4	4, 8, 16 , 32	30	-	-
<i>Structure prediction</i>						
ESMFold + TTT	4e-4	4	4, 8, 32, 64	30 (max pLDDT)	4, 8 , 32	8, 16, 32
ESM3 + TTT	1e-4, 4e-4, 1e-3	2	1 , 4, 16	50 (max pLDDT)	-	-
<i>Function prediction</i>						
TerpeneMiner + TTT	4e-4 , 1e-3	2	2, 4, 8	30	-	-
Light attention + TTT	4e-4, 1e-3, 3e-3	2	2, 4	30	-	-

B.1 DETAILED TEST PERFORMANCE

In this section, we provide details on the test performance. Specifically, Table 5 shows that test-time training (TTT) 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 an example illustrating how TTT substantially improves the correlation between ESM2-predicted fitness and ground-truth stability by better identifying disruptive mutations in the protein core (Figure 7).

Next, Figure 9 shows the distribution of TTT effects: in many cases, TTT has minimal impact on performance; often, it leads to substantial improvements; and in rare cases TTT results in a decrease in performance. This positions TTT as a method for enhancing prediction accuracy, while a comprehensive analysis of its failure modes remains an important direction for future research. While we demonstrate these effects using a protein folding example, we observe a similar distribution of TTT impact across the tasks.

We also observe that the overall trend of TTT generally leads to improved performance, with robust consistency across random seeds. However, the progression of the performance curve can be rugged, particularly in classification tasks, where substantial changes in the underlying representations are required to shift the top-predicted class in the discrete probability distribution (Figure 11).

B.2 VALIDATION PERFORMANCE

This section discusses the performance of test-time training (TTT) on validation data. Table 6 illustrates the validation performance of all tested methods for fitness prediction on our newly constructed MaveDB dataset. TTT enhances the performance of all the methods.

The primary focus of the section is hyperparameter tuning. Table 4 provides the grid of hyperparameters explored for each model and its size. Figure 12 demonstrates the trend of hyperparameter tuning with optimal hyperparameter combination balancing underfitting and overfitting to a single test protein. While most hyperparameter configurations lead to overall improvements when using TTT, poorly chosen hyperparameters can have detrimental effects due to rapid overfitting. However, with a reliable predicted confidence measure, such as pLDDT, the appropriate TTT step can be selected to mitigate overfitting. Figure 13 demonstrates that when using ESM3 + TTT with pLDDT-based step selection for protein folding, all hyperparameter configurations result in improved performance compared to the base ESM3 model.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

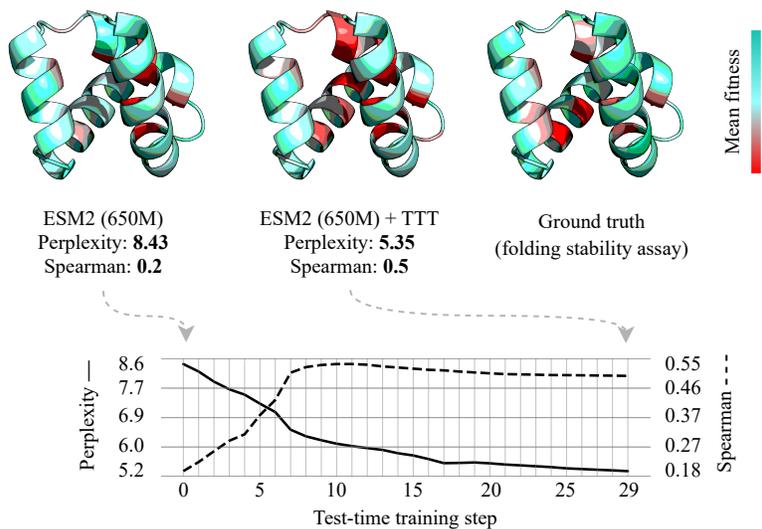


Figure 7: **Example of test-time training (TTT) applied to fitness prediction.** Fitness predictions from ESM2 (650M) show poor correlation with experimental fitness values in the ProteinGym test set measured by the stability assay “UBR5_HUMAN_Tsuboyama_2023_1I2T” (Tsuboyama et al., 2023) (left). ESM2 + TTT 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 TTT-enhanced 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.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

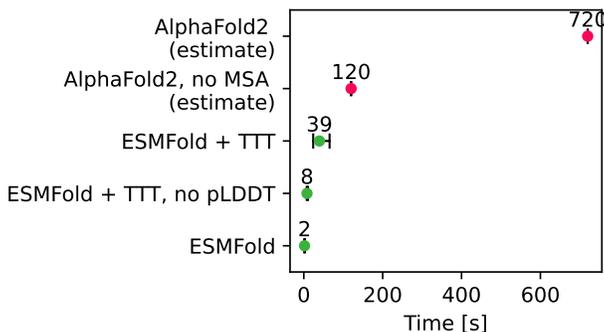
Table 5: **Test-time training (TTT) performance on ProteinGym depending on MSA depth.** MSA depth reflects the number of available proteins similar to the target protein and, when using large protein language models, can be interpreted as a measure of the representation of similar proteins in the training data (Appendix A.1.2). TTT 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) + TTT (Ours)	0.3407 \pm 0.00014	0.2445	0.3144	0.4598
SaProt (35M) (Su et al., 2023)	0.4062	0.3234	0.3921	0.5057
SaProt (35M) + TTT (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) + TTT (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) + TTT (Ours)	0.4583 \pm 0.00001	0.3954	0.4501	0.5439

Table 6: **Performance of test-time training (TTT) on the MaveDB dataset.** In this work, we use our newly constructed MaveDB benchmark as a validation fold for tuning the hyper-parameters of TTT for fitness prediction. For computational efficiency, we only select a subset of 50 proteins (Appendix A.1.1) and do not run TTT 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) + TTT (Ours)	0.4593
ESM2 (650M) (Lin et al., 2023)	0.4568
ESM2 (650M) + TTT (Ours)	0.4604
SaProt (650M) (Su et al., 2023)	0.4926
SaProt (650M) + TTT (Ours)	0.4926
SaProt (35M) (Su et al., 2023)	0.5251
SaProt (35M) + TTT (Ours)	0.5271

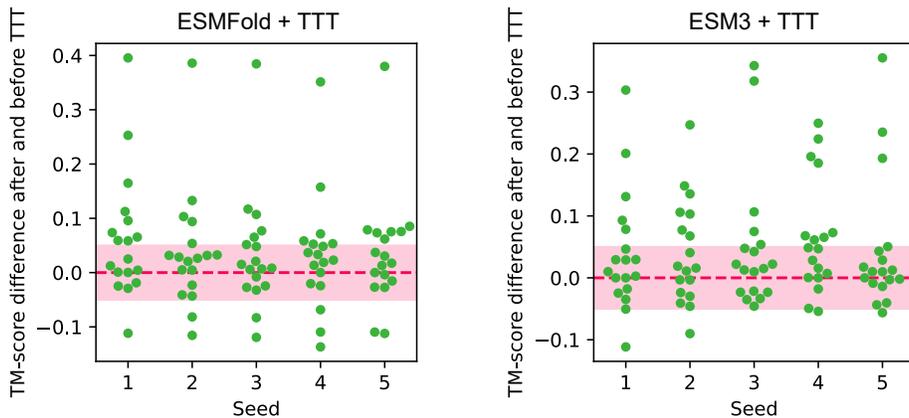
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471



1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482

Figure 8: Running time of ESMFold + TTT. 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: $2 \times 60 = 120$ seconds), with additional MSA construction taking at least 10 minutes using standard pipelines (AlphaFold2: $2 \times 60 + 10 \times 60 = 720$ seconds). ESMFold + TTT (30 steps) involves test-time training parameter updates with LoRA, along with forward passes at each TTT step to estimate pLDDT and select the structure with the highest predicted confidence. Disabling pLDDT significantly reduces computational overhead (ESMFold + TTT, no pLDDT compared to ESMFold + TTT), but may require careful parameter tuning (Appendix B.2). Overall, ESMFold + TTT maintains the speed advantage of ESMFold, and is significantly faster than AlphaFold2.

1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501



1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

Figure 9: Per-protein performance of ESMFold + TTT and ESM3 + TTT on the CAMEO test set. The y-axis shows the change in TM-score after applying test-time training (TTT), 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 TTT 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 TTT to ESM3 results in 6 improvements out of 18 proteins, with 1 case of degradation.

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

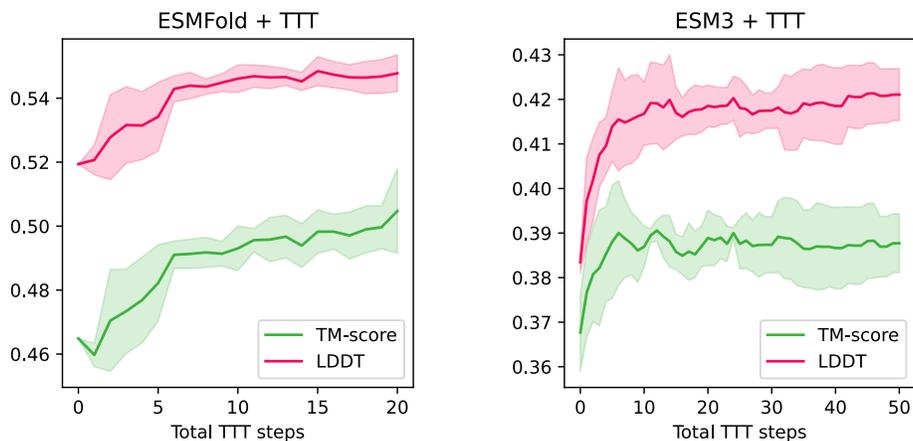


Figure 10: **Test performance of ESMFold + TTT and ESM3 + TTT on the CAMEO test set depending on the total number of TTT steps.** The x-axis shows the averaged performance across all test proteins, with error bars representing the standard deviation across five random seeds. The y-axis metrics correspond to the structure with the highest pLDDT score up to the given step. While an increased number of TTT steps generally enhances performance, only a few TTT steps (e.g., five) may suffice to achieve significant performance improvement.

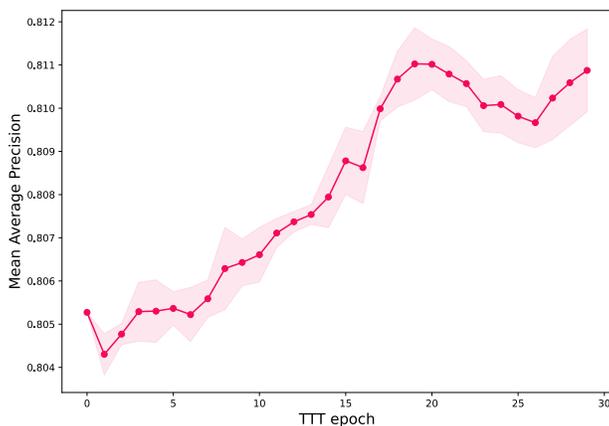


Figure 11: **Test performance of TerpeneMiner + TTT across fine-tuning steps.** The performance is averaged across all 512 proteins in the dataset, with error bars representing the standard deviation across 5 random seeds.

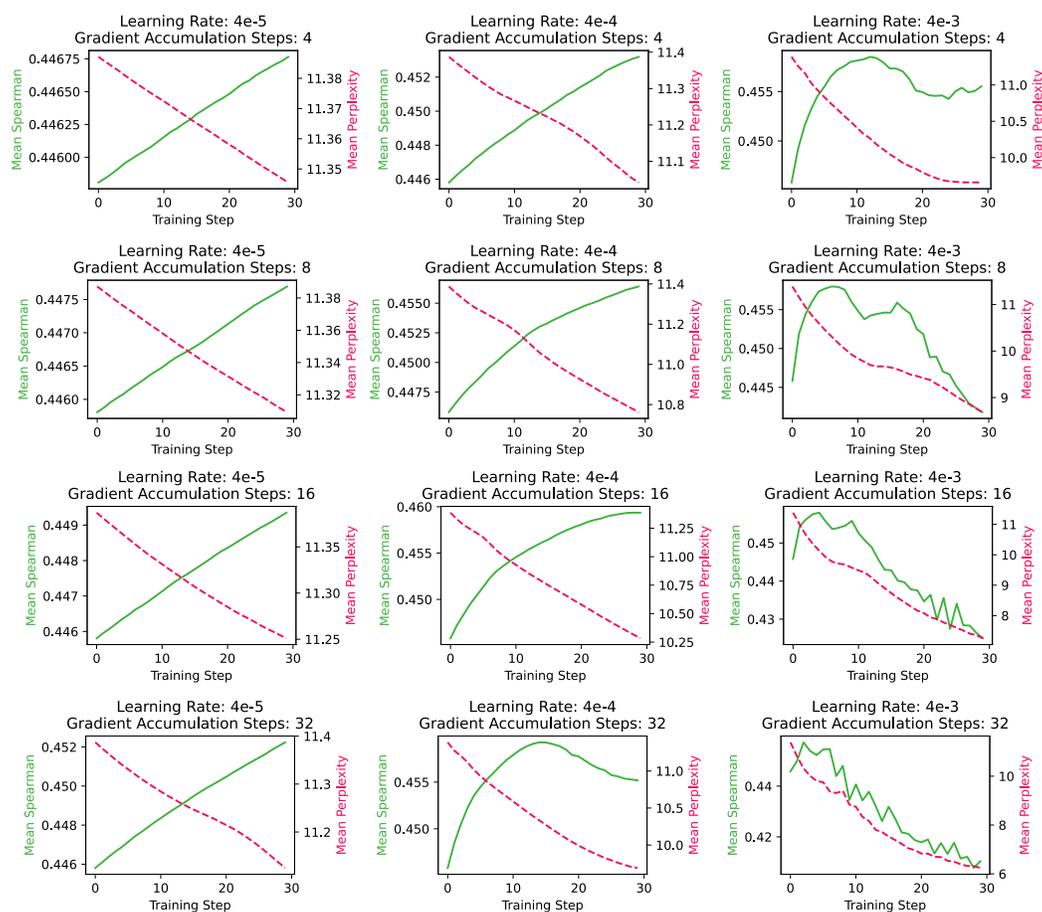


Figure 12: **Dependence on hyperparameters in test-time training for fitness prediction.** Each plot shows the progression of Spearman correlation (green) increasing alongside a decrease in perplexity (pink) for each TTT step, averaged across all assays in the MaveDB validation dataset. The model used is ESM2 (35M) + TTT, 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 test 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 test protein. **Notably, the figure demonstrates that, although TTT involves three main hyperparameters (batch size, learning rate, and the number of TTT 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.**

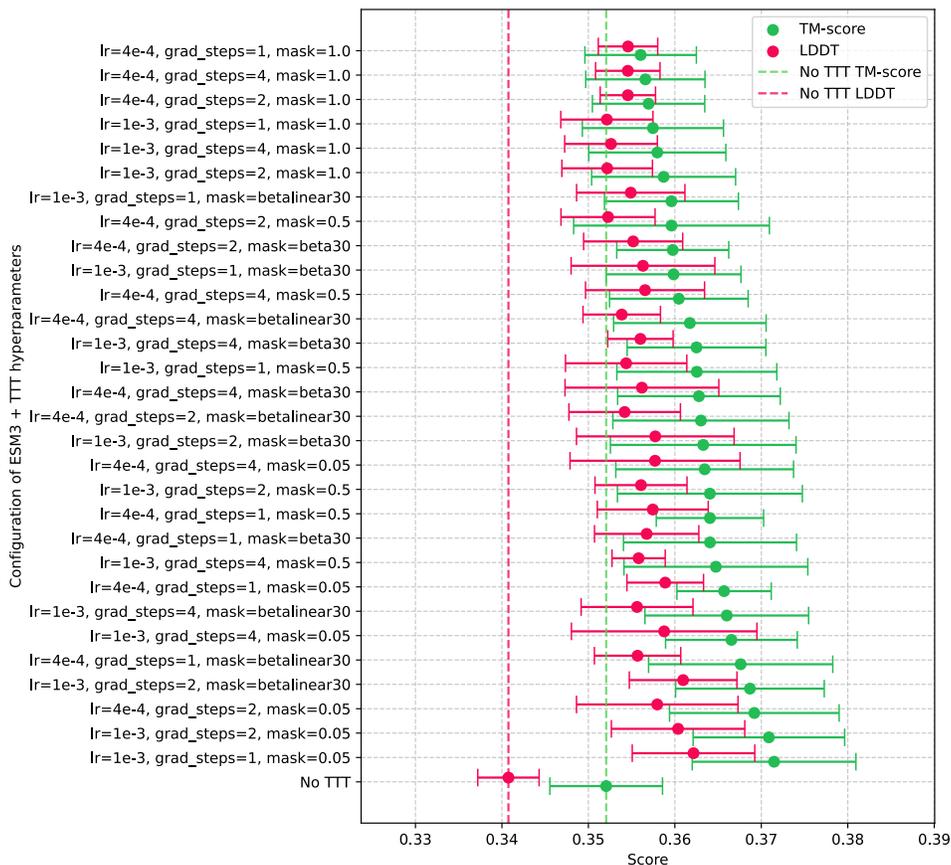


Figure 13: **Hyperparameter search for protein structure prediction with ESM3 + TTT.** We conducted a comprehensive grid search based on three key hyperparameters: learning rate (denoted as “lr”), number of gradient accumulation steps (denoted as “grad_steps”; with the batch size of two), and masking 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 TTT”, shows the performance of ESM3 without TTT. The results indicate that ESM3 + TTT 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 TTT”) to compute the results on the test fold. For the hyperparameter search, we used 30 TTT steps instead of 50 to reduce computation time.