

1 Understanding Language Model Scaling on Protein Fitness Prediction

2 Chao Hou^{1,#}, Di Liu², Aziz Zafar², Yufeng Shen^{1,2,3,4,#}

3 1 Department of Systems Biology, Columbia University Irving Medical Center, New York, NY 10032

4 2 Department of Biomedical Informatics, Columbia University Irving Medical Center, New York, NY 10032

5 3 JP Sulzberger Columbia Genome Center, Columbia University, New York, NY 10032

6 4 Program for Mathematical Genomics, Columbia University Irving Medical Center, NY 10032

7 # Corresponding author: ch3849@cumc.columbia.edu (C.H.), ys2411@cumc.columbia.edu (Y.S.)

8

9 Abstract

10 Protein language models, and models that incorporate structure or homologous sequences, estimate
11 sequence likelihood $p(sequence)$ that reflects the protein fitness landscape and is commonly used in
12 mutation effect prediction and protein design. It is widely believed in deep learning field that larger model
13 performs better across tasks. However, for fitness prediction, language model performance declines
14 beyond a certain size, raising concerns about their scalability. Here, we showed that model size, training
15 dataset, and stochastic elements can bias the predicted $p(sequence)$ away from real fitness. Model
16 performance on fitness prediction depends on how well $p(sequence)$ matches evolutionary patterns in
17 homologs, which is best achieved at a moderate $p(sequence)$ level for most proteins. At extreme predicted
18 wild-type sequence likelihoods, models predict uniformly low or high likelihoods for nearly all mutations,
19 failing to reflect the real fitness landscape. Notably, larger models tend to predict proteins with higher
20 $p(sequence)$, which may exceed the moderate range and thus reduce performance. Our findings clarify
21 the scaling behavior of protein models on fitness prediction and provide practical guidelines for their
22 application and future development.

23 **Keywords:** protein fitness landscape, protein language model, self-supervised learning, sequence
24 likelihood, mutation effect.

25

26 Introduction

27 Characterizing the protein fitness landscape—how mutations affect protein function, abundance, activity,
28 and interaction—is a central challenge in biology. It is crucial for elucidating the mechanisms underlying
29 diseases, advancing precision medicine, guiding viral surveillance, and advancing protein design and
30 engineering. While deep mutational scanning¹ (DMS) experiment has been applied to measure mutation
31 effects on diverse proteins^{2,3}, it's time-consuming, labor-intensive, and limited to molecular effects that
32 are easy to assay. To complement experimental efforts, supervised machine learning methods have been
33 developed by training on curated mutation datasets⁴⁻⁶. However, these datasets are limited in size and
34 biased toward functionally important genes⁴, restricting the models' generalizability and robustness⁷. As
35 a result, there is an urgent need for predictive models capable of estimating the protein fitness landscape
36 without training on curated mutation datasets.

37 In recent years, self-supervised models have been developed for zero-shot fitness prediction by estimating
38 sequence likelihood $p(sequence)$ (Throughout this manuscript, “sequence likelihood” and $p(sequence)$
39 refer to the likelihood of wild-type sequence unless otherwise noted for mutant sequences.): the probability
40 of a protein sequence under the learned distribution of natural proteins. Self-supervised models often
41 match or even surpass the performance of supervised models^{8,9}. Representative models include protein
42 language models (pLMs) such as ESM2¹⁰; multi-sequence alignment (MSA)-based models like MSA-
43 Transformer¹¹ and EVE⁸; inverse folding models such as ESM-IF1¹², which predict sequence from structure;
44 and hybrid models like ESM3¹³, which integrate sequence, structure, and other information. Based on their
45 training data, these methods fall into two main categories: general models trained on massive datasets of

1 tens to hundreds of millions of proteins from diverse protein families, and family-specific models trained
2 on MSA of individual protein family (e.g., EVE, which requires training a separate model for each MSA).
3 MSA-Transformer integrates both family-specific information and general information across diverse
4 protein families by being trained on millions of MSAs¹¹.

5 These models are trained to maximize the likelihood of training protein sequences using strategies such
6 as masked or next token prediction¹⁴, conditioned on sequence, structure, or MSA (**Figure 1A**). While
7 some models incorporate additional training objectives¹³, only sequence prediction is used to predict
8 fitness (see Methods for details)³. Fitness of a mutation is estimated by comparing the predicted likelihoods
9 of the mutant and wild-type sequences^{3,9} (**Figure 1A**), usually by calculating the log-likelihood ratio (LLR).
10 In this framework, an LLR close to zero means that the mutation is nearly as fit as the wild-type, implying
11 a neutral mutation effect, whereas a strongly negative LLR indicates the mutation is much less fit and
12 potentially deleterious. For models trained via masked prediction, computing the $p(sequence)$ is intractable.
13 Instead, pseudo-likelihood and LLR calculated from marginal approaches are used (see Methods for
14 details)^{3,9}.

15 Among these models, pLMs gain particular attention for their strong performance and minimal input
16 requirement—they only require sequence as input, without needing MSA or structure. In deep learning
17 field, scaling up models is a common strategy to improve performance on downstream tasks¹⁵. pLM
18 scaling has proven effective for masked or next residue prediction and structure modeling^{10,16}; however,
19 for fitness prediction, model performance declines beyond a certain size^{3,17}: ESM2-650M (million
20 parameters) and 3B (billion parameters) outperform the 15B model (**Table S1**), xTrimoPGLM-3B¹⁸
21 outperforms the 10B and 100B models, and autoregressive pLMs ProGen3-1B and 3B outperform the
22 10B and 46B models¹⁶. This raises a key question in model development: why does scaling up pLMs not
23 consistently improve performance on fitness prediction? While pLMs are typically released in multiple
24 sizes, enabling direct investigation of the scaling behavior, other models are usually released in one size,
25 making it unclear whether the scaling trend applies to them. Beyond scaling, several important questions
26 remain: Under what conditions are these models most effective? How to choose appropriate models for
27 a given protein?

28 In this study, we systematically investigated the relationship between sequence likelihood and the
29 performance of fitness prediction across diverse models. We found that various factors unrelated to
30 fitness—including model size, training dataset, and stochastic elements—can influence general model-
31 predicted likelihoods, making them uninformative for fitness prediction in extreme cases. By analyzing
32 large fitness benchmarks^{2,3}, we found that the performance of general models depends on how well their
33 predicted $p(sequence)$ aligns with evolutionary patterns in homologs. Notably, they perform best at a
34 moderate $p(sequence)$ level for most proteins, which explain the scaling behavior of pLMs as medium-
35 sized models like ESM2-650M predict more proteins with moderate $p(sequence)$. Our findings clarify how
36 likelihood-based self-supervised models predict fitness and lay the groundwork for developing next-
37 generation fitness predictors.

38

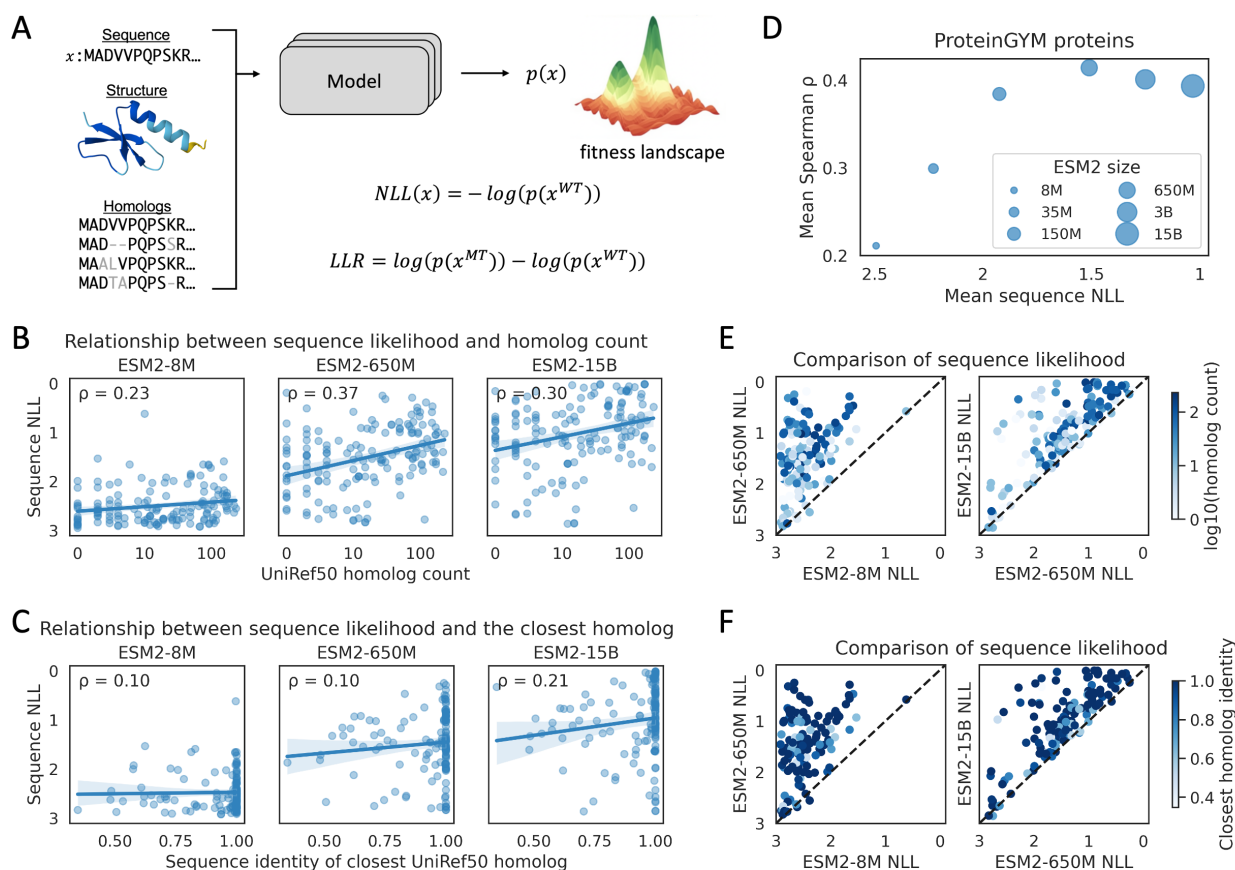
39 **Results**

40 **General model-predicted $p(sequence)$ is influenced by factors unrelated to protein fitness.**

41 As general models trained on large datasets function as black boxes, we examined whether factors
42 unrelated to fitness can influence the predicted $p(sequence)$. We first investigated two key factors in
43 scaling up models: the size of the training dataset and the number of trainable parameters. To do this, we
44 analyzed predictions from several models on 154 fitness measurements from deep mutational scanning
45 (DMS) experiments in the ProteinGYM³ benchmark (see Methods for data filtering), with a primary focus
46 on six ESM2 models spanning 8 million to 15 billion parameters. ProteinGYM includes proteins from
47 diverse taxa and five functional categories, spanning mutation effects on activity, binding, stability,
48 organismal fitness, and expression.

49 Larger training datasets include more protein families and more homologs within each family. Since
50 dissimilar proteins have little influence on the predicted likelihood of a protein¹⁹, we focused on the role of

1 its homologs. We identified UniRef50²⁰ (the ESM2 training dataset, version2021_04) homologs with over
 2 20% sequence identity and 80% coverage to proteins in the ProteinGYM benchmark. By analyzing the
 3 number of homologs, we found that proteins with more homologs tend to have higher predicted likelihoods
 4 (**Figure 1B**, quantified using negative log-likelihood (NLL, $-\log p(\text{sequence})$)). We note that this trend is not
 5 indirectly driven by conservation level differences of protein families (**Figure S1A**). However, proteins with
 6 high predicted likelihoods do not necessarily have many UniRef50 homologs (**Figure 1B**). We then
 7 analyzed the sequence identity of the closest homologs and observed that proteins with high predicted
 8 likelihoods often have highly similar homologs in UniRef50. But the presence of such homologs does not
 9 guarantee high predicted likelihoods (**Figure 1C**). We also analyzed the summed sequence identity \times
 10 coverage of all homologs, which provides a more quantitative measure of the overall similarity of homologs
 11 in the training set. The results are similar to those obtained using homolog count (**Figure S1B**). Additionally,
 12 we observed similar relationships between homologs in the training set and sequence likelihoods for
 13 SaProt²¹, a structure-informed model (**Figure S2A-B**). Overall, homologs in the training set can influence
 14 model-predicted likelihoods, but the relationship is complex.
 15



16

17 **Figure 1. Model-predict sequence likelihood is influenced by factors unrelated to protein fitness.**

18 **A**, Overview of the calculation of sequence likelihood and log-likelihood ratio (LLR). Models are trained to predict sequence likelihood
 19 $p(x)$ using information from sequence, structure, and homologs, negative log-likelihood (NLL) is usually used as the training loss. The
 20 lower-left panel shows a multiple sequence alignment (MSA), the dashed lines indicate alignment gaps, and the grey-shaded amino
 21 acids denote positions in homologous sequences that differ from the query sequence (the first sequence in the MSA). For mask-
 22 prediction-based models, pseudo-likelihood is used. LLR is calculated based on the entire sequence for generative models, and
 23 masked residues for mask-prediction-based models. **B**, Relationship between predicted sequence likelihood and the number of
 24 homologs. Each point represents a protein; the x-axes show the number of UniRef50 homologs in the log scale. Homologs are
 25 defined as those with $\geq 20\%$ sequence identity and $\geq 80\%$ coverage. The curves represent linear regressions, with the shaded areas
 26 indicating the 95% confidence intervals. ρ represents Spearman correlation. **C**, Relationship between predicted sequence likelihood
 27 and sequence identity of the closest homolog. **D**, The weighted mean performance on fitness prediction, and mean predicted

1 sequence likelihood of ESM2 models on proteins in 154 ProteinGYM experiments. Point size indicates ESM2 model size. The
2 performance is first averaged within each of the five functional categories in ProteinGYM, and the weighted mean performance is
3 computed as the average of five category-specific performances. These values are reported in Table S1. **E-F**, Comparison of
4 sequence likelihoods predicted by different ESM2 models. Each point represents a protein; colors indicate the log-scaled number
5 of homologs (**E**) and sequence identity of the closest homolog (**F**).

6
7 We then analyzed model size. These models are typically trained using the NLL of masked or next token
8 as the loss function. As larger models achieve lower training loss, they tend to predict higher sequence
9 likelihoods^{10,16,18}. We observed this trend for the proteins in the ProteinGYM benchmark (**Figure 1D**).
10 Furthermore, we examined the magnitude of likelihood increase as model size scaled up and observed
11 substantial variability: some proteins showed little or no increase, while others exhibited large gains
12 (**Figure 1E-F**). Notably, the magnitude of likelihood increase between larger and smaller models is not
13 clearly associated with the likelihoods from the smaller model, nor with the number or similarity of
14 homologs in the training set (**Figure 1E-F**). We observed similar results for SaProt, xTrimoPGLM¹⁸, and
15 ProGen3¹⁶ (**Figure S2C-D**).

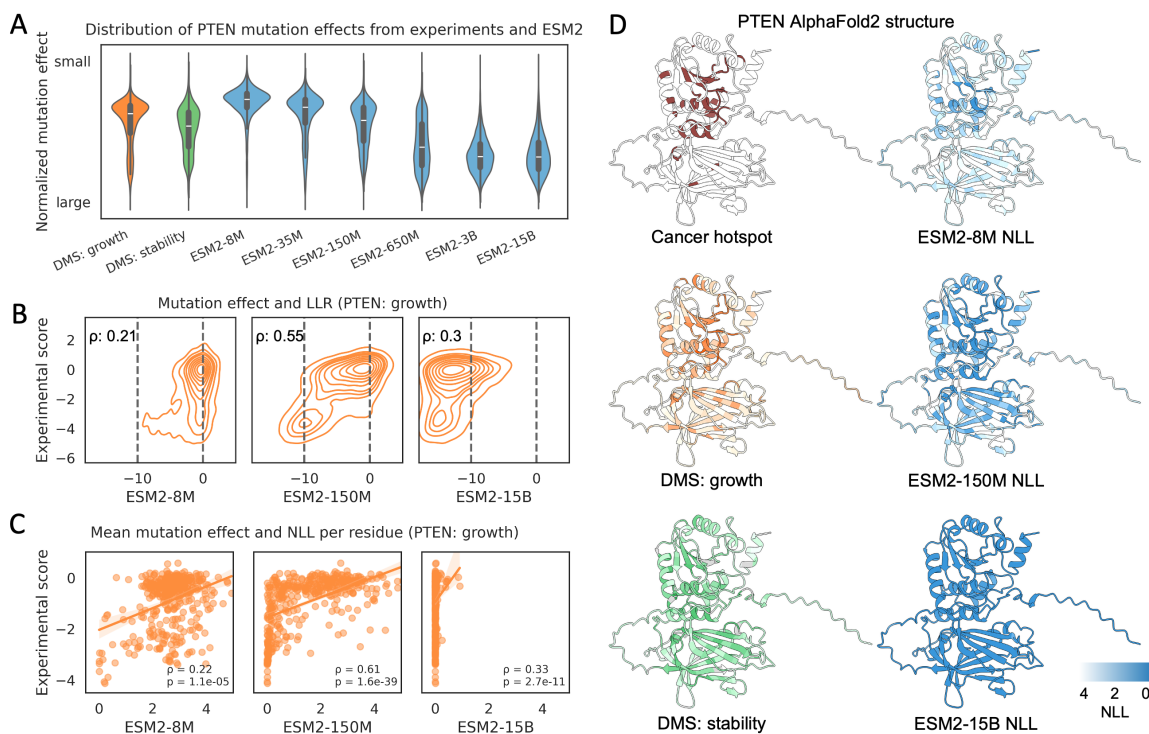
16 These results suggest that, beyond the homologs in training dataset and model size, additional factors
17 influence model-predicted likelihoods. One such factor could be the stochastic elements introduced
18 during model training, including parameter initialization, data shuffling, and masked residue sampling. To
19 directly assess the impact of these stochastic elements, we analyzed five ESM1v²² models, which share
20 the same architecture and training dataset but were trained with different random seeds. We found that
21 approximately 10% of ProteinGYM proteins show NLL differences greater than 0.5 among five ESM1v
22 models (**Figure S3**), and the maximum observed difference is 1.2: the protein Mafg (UniProt ID: O54790)
23 has a predicted sequence NLL of 1.1 in ESM1v_5 but 2.3 in ESM1v_2. This result suggests that stochastic
24 elements can lead models to converge on different local minima, resulting in variability in predicted
25 sequence likelihoods for some proteins. Overall, model size, training dataset, stochastic elements, and
26 other unknown factors complicate the interpretation of general model-predicted likelihoods.

27
28 **Magnitude of predicted $p(\text{sequence})$ affects fitness estimation by influencing LLR values.**

29 As model-predicted sequence likelihood is affected by many factors, the LLR value, directly tied to
30 sequence likelihood (LLR is the NLL difference between wild-type and mutant sequences, **Figure 1A**) and
31 used to estimate fitness, is also affected. This may impact fitness prediction performance. To investigate
32 this, we focused on model size, as other factors are difficult to control given models we have access to.
33 Larger models, achieving lower training loss, tend to assign higher probability to the wild-type amino acid
34 and lower probabilities to the others (since the total probability per site sum to one), resulting in LLRs with
35 larger magnitudes. In extreme cases, a non-informative model that predicts equal probability to all 20
36 amino acids yields LLRs of zero for all mutations, while an overconfident model that assigns a probability
37 of one to the wild-type and zero to all others produces LLRs of negative infinity. Although certain proteins
38 may predominantly harbor neutral or deleterious mutations, the collapsed LLR distributions in these two
39 extreme scenarios are uninformative for most proteins.

40 To illustrate this, we examined predictions of six ESM2 models on PTEN (phosphatase and tensin
41 homolog), one of the most extensively studied proteins, for which DMS experiments of both cell growth²³
42 and protein stability²⁴ are available. Across both DMS experiments, the distribution of mutation effects is
43 clearly bimodal, with approximately 20% of mutations exhibiting deleterious effects (**Figure 2A**). However,
44 the distributions of ESM2 predicted LLR vary substantially with model size. Smaller models predict LLRs
45 clustered near zero, while larger models predict strongly negative LLRs to most mutations (**Figure 2B**,
46 **S4A**), mirroring the two extreme scenarios we described above. Notably, the medium-sized model ESM2-
47 150M reproduces the bimodal distribution observed in experiments (**Figure 2A, S5A**) and achieves the
48 best performance (**Figure S5B**), with Spearman correlations of 0.55 for growth and 0.46 for stability. In
49 contrast, both the smallest (ESM2-8M) and largest (ESM2-15B) models yield correlations below 0.3 in both
50 DMS experiments (**Figure 2B, S4A, S5B**). We also observed that xTrimoPGLM-1B model best captures
51 the distribution of experimental mutation effects and outperforms larger models with 3B to 100B

1 parameters (**Figure S5C-D**). We note that the smallest available xTrimoPGLM model is 1B, and thus
 2 comparisons to smaller models are not possible. We further investigated this relationship at the residue
 3 level by comparing the mean mutation effect per residue with the predicted probability of the wild-type
 4 residue. Residues assigned high probability by the model tend to have strongly negative LLRs for
 5 mutations, indicating that such residues are mutation sensitive. ESM2-15B assigns high probability
 6 approaching one to nearly all residues, while ESM2-8M assigns high probability to very few residues—
 7 both failing to capture the experimentally observed distribution of mutation-sensitive residues (**Figure 2C-
 8 D, S4B**). In contrast, ESM2-150M, predicts residue probability that better reflect mean mutation effects
 9 (**Figure 2C-D, S4B**). Additionally, we analyzed PTEN cancer hotspot mutations²⁵, and found that predicted
 10 residue likelihood from medium-sized ESM2 models (35M and 150M) outperform other ESM2 models at
 11 identifying residues with hotspot mutations (**Figure S5E**). These results suggest that medium-sized ESM2
 12 models more accurately capture the residue-level importance in PTEN. We observed similar results in
 13 other proteins that also exhibit a rise-then-fall performance trend with increasing ESM2 model size (**Figure
 14 S4C-H**).
 15



16
 17 **Figure 2. Distributions of experimental and predicted PTEN mutation fitness.**

18 **A**, Distributions of normalized mutation effects from DMS experiments and ESM2 predictions. Experimental effects and ESM2 LLRs
 19 are normalized to the range of 0–1 for visualization. ESM2 LLRs are calculated using the masked marginal approach. **B**, Relationship
 20 between ESM2-predicted LLRs and experimental effects for PTEN mutations. The x-axis represents predicted LLR, the y-axis
 21 represents log-scaled, wild-type-normalized fitness measurements based on the growth of humanized yeast. ρ : Spearman
 22 correlation. **C**, Relationship between ESM2-predicted probability per residue (quantified using NLL) and mean experimental effects
 23 for mutations at each residue (residues with at least 10 mutations are shown). The x-axis represents predicted NLL per residue, the
 24 y-axis represents mean experimental scores per residue from ProteinGYM, each point represents a residue. Spearman correlation
 25 and corresponding p-values are shown. **D**, Visualization of mutation sensitivity and ESM2-predicted NLL per residue on the
 26 AlphaFold2-predicted PTEN structure. For two DMS experiments, the mean mutation effect at each residue is shown (grey color
 27 indicates no mutation at the site). Darker colors indicate stronger mutation effects or higher ESM2 predicted likelihood (i.e., lower
 28 NLL).

29
 30 **General model performance on fitness prediction peaks at a moderate level of $p(\text{sequence})$.**

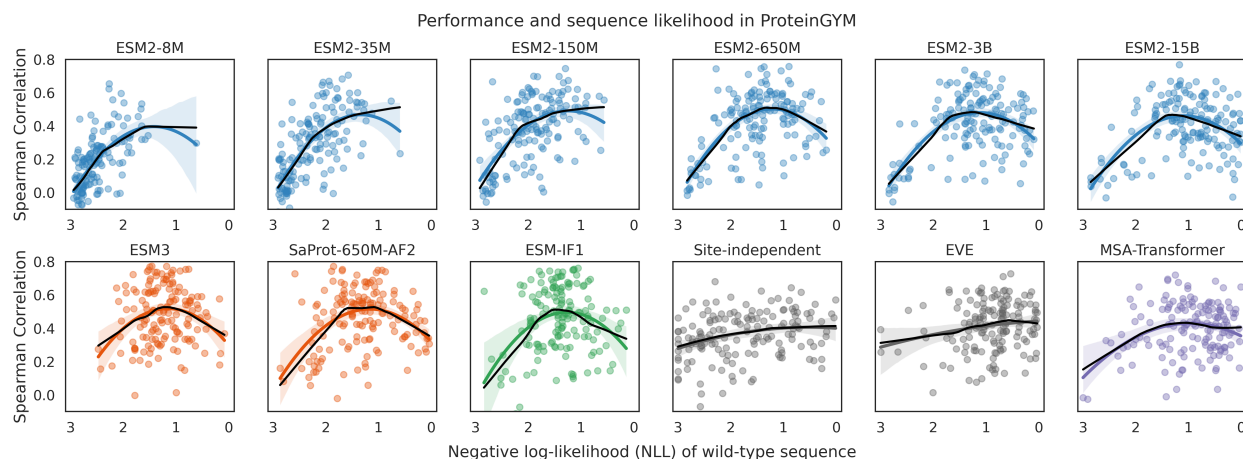
1 The above analyses highlight a critical caveat when using general models for fitness prediction: their
2 predicted likelihoods can be influenced by unrelated factors and may not reliably reflect real fitness. To
3 systematically investigate the relationship between model-predicted sequence likelihood and performance
4 on fitness prediction, we evaluated models on 154 DMS experiments from the ProteinGYM benchmark³
5 (see Methods for dataset filtering).

6 We first evaluated six ESM2 models. Although larger models consistently predict higher sequence
7 likelihood (**Figure 1D**), higher sequence likelihood does not always lead to better performance. Notably,
8 the ESM2-650M model outperforms the larger 3B and 15B models (**Figure 1D, Table S1**). By comparing
9 wild-type sequence likelihoods with performances of 154 experiments, we observed a bell-shaped
10 relationship using the non-parametric LOWESS²⁶ modeling (locally weighted scatterplot smoothing, black
11 curves in **Figure 3**), with performance peaking at a moderate likelihood level—corresponding to a wild-
12 type sequence NLL of approximately 1.2 (**Figure 3**, $p(\text{sequence}) \approx 0.3$). Because second-order polynomial
13 regressions closely match the non-parametric LOWESS trends (**Figure 3**), we also used it to describe the
14 bell-shaped relationship and reported its 95% confidence intervals. Notably, within the optimal likelihood
15 range, all ESM2 models, regardless of size, perform comparably (**Figure 3**). We observed the same trend
16 for ESM2 models on the mega-scale protein folding stability dataset² (**Figure S6A**), indicating that this
17 relationship is not specific to ProteinGYM. These results indicate that the level of sequence likelihood,
18 rather than model size, is the primary determinant of model performance on fitness prediction. While more
19 homologs in the ESM2 training set can lead to higher predicted sequence likelihoods (**Figure 1B**), we
20 found that the performance drop at high likelihood is not caused by homolog overrepresentation in the
21 training set (**Figure S1C**, see **Figure S2F** for SaProt).

22 We further evaluated a broad range of general models that predict fitness using sequence likelihood. These
23 included Transformer encoder-based mask language models ESM1v²² and ESMC²⁷, whose training
24 datasets are different from that of ESM2; the convolution-based model CARP²⁸; and Transformer decoder-
25 based generative models ProGen3¹⁶ and RITA²⁹. We also evaluated structure-sequence-hybrid models
26 ESM3¹³, ProSST³⁰, and SaProt²¹, as well as the inverse folding model ESM-IF1¹² (see Methods for details).
27 Despite differences in architecture, input modalities, and training strategies, all these general models
28 exhibit the bell-shaped relationship between performance and wild-type sequence likelihood (**Figure 3**,
29 **S7A**; see Discussion for decoder-based pLMs). Remarkably, peak performances are achieved at the
30 similar level of sequence likelihood (**Figure 3, S7A**). For ProSST and small ESM2 models, only one side of
31 the bell-shaped trend is observed, as almost no proteins exhibit predicted sequence likelihoods beyond
32 the optimal range due to limited model capacity. Because ProSST shares an almost identical model design
33 with SaProt, we expect it to exhibit the same bell-shaped relationship. In addition, for all LLR-based
34 models, fitness predictive power necessarily vanishes when the wild-type sequence NLL is zero: in this
35 limit, models predict LLRs of negative infinity to all mutations. Consequently, the performance curve must
36 pass through the origin (NLL = 0, fitness prediction Spearman correlation = 0). In practice, however,
37 sequence NLL values of zero are not observed for these models. We used the masked marginal approach
38 for models trained via masked token prediction. We note that we also observed the bell-shaped trend
39 when these models are applied using the wild-type marginal approach (**Figure S7B-D**, see Methods for
40 details).

41 Beyond general models trained on diverse protein families, we also analyzed family-specific models. These
42 included simple frequency-based site-independent models that treat each position independently, with or
43 without homologous sequence weighting (see Methods for details), and the variational autoencoder-based
44 model EVE⁸, which captures inter-residue dependencies. These models are trained independently on MSA
45 of each protein family with no or less parameters, making them less susceptible to the unrelated factors
46 that affect general models. Notably, these family-specific models do not exhibit the bell-shaped
47 relationship between performance and sequence likelihood (**Figure 3, S7A**). For MSA-Transformer¹¹ which
48 integrates both general and family-specific information, the bell-shaped trend is present but less
49 pronounced (**Figure 3**).

1

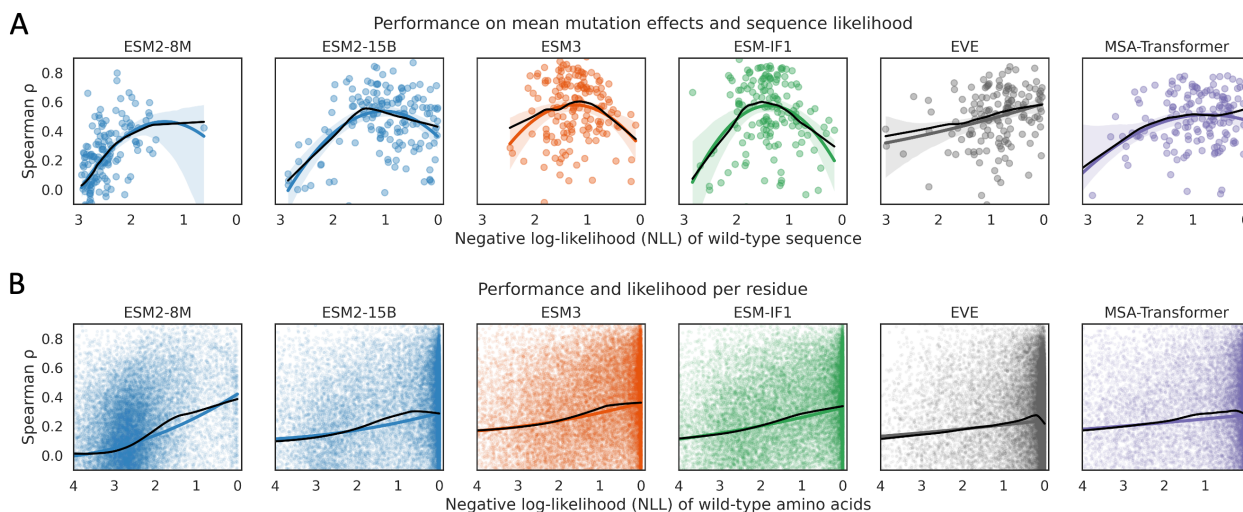


2

3 **Figure 3. Relationship between fitness prediction performance and model-predicted wild-type sequence likelihood.**

4 The y-axes show the Spearman correlation between model predicted LLRs and experimental effects, the x-axes represent the
 5 negative log-likelihood (NLL) of wild-type sequences. LLRs of ESM2, ESM3, SaProt, and MSA-Transformer are calculated using the
 6 masked marginal approach. LLRs of ESM-IF1 and EVE are calculated using the full sequence (see Methods for details). The LLR of
 7 the site-independent model is computed as the log difference in amino acid frequency at each site. Only residues with mutations are
 8 included in the NLL calculation. Each point represents an experiment from ProteinGYM (154 experiments after filtering). The black
 9 curves show LOWESS modeling, while the colored curves show second-order polynomial regressions, with shaded areas indicating
 10 the 95% confidence intervals, all regression analyses were performed using the Python package Seaborn. Colors indicate model
 11 types: blue: protein language models; orange: hybrid models, green: inverse-folding models, grey: family-specific models trained on
 12 MSA of individual protein family, and purple: MSA-Transformer that integrate both general and family-specific information.

13



14

15 **Figure 4. Model performance on mean mutation effects and mutation effects per residue.**

16 **A.** The y-axes show the Spearman correlation between mean LLRs and mean experimental mutation effects per residue in each
 17 protein, reflecting model understanding of protein context. The x-axes represent the NLL of wild-type sequences. Only residues with
 18 mutations are included in the NLL calculation. Each point represents a ProteinGYM experiment. **B.** The y-axis show the Spearman
 19 correlation between LLRs and experimental effects of all 19 mutations per residue, reflecting model understanding of substitution
 20 specificity. The x-axes represent the negative log predicted probability of each residue. Each point represents a residue. The black
 21 curves show LOWESS modeling, while the colored curves show second-order polynomial regressions, with shaded areas
 22 representing the 95% confidence intervals.

23

24 **The bell-shaped relationship arises from models' varying ability to capture context information.**

1 Fitness prediction requires models to capture both context information (i.e., the sequence and structural
2 context that determine mutation sensitivity) and substitution specificity (i.e., how well different amino acids
3 fit a given context). To investigate the origin of the bell-shaped relationship between performance and
4 sequence likelihood, we disentangled two components. Context understanding is quantified by comparing
5 the mean LLR and the mean experimental effect for mutations on each residue within a protein.
6 Substitution specificity is assessed by correlating LLRs with experimental effects of 19 mutations at each
7 residue.

8 By comparing model performance on the two components with predicted likelihood, we found that general
9 models exhibit the bell-shaped relationship between context understanding (mean mutation effect
10 prediction performance) and wild-type sequence likelihood, with performance peaking at a similar
11 likelihood range (**Figure 4A, S8A, S6B**). Extreme cases are exemplified by our results for PTEN and other
12 proteins (**Figure 2C, S4**). In contrast, per-residue performance increases monotonically with predicted
13 likelihood (i.e., the predicted probability of each residue; **Figure 4B, S8B, S6C**; see Discussion for
14 explanation). Family-specific models, however, do not display the bell-shaped trend for either component
15 (**Figure 4, S8**). Notably, all models show substantially stronger performance in understanding context than
16 substitution specificity on 154 experiments (**Figure 4, S8**). None of the models evaluated in ProteinGYM
17 achieves mean Spearman correlation above 0.3 for substitution specificity. These findings indicate that
18 the fitness predictive power of current models primarily stems from their ability to capture protein context,
19 which also underlies the bell-shaped relationship observed between fitness prediction performance and
20 sequence likelihood in general models.

21

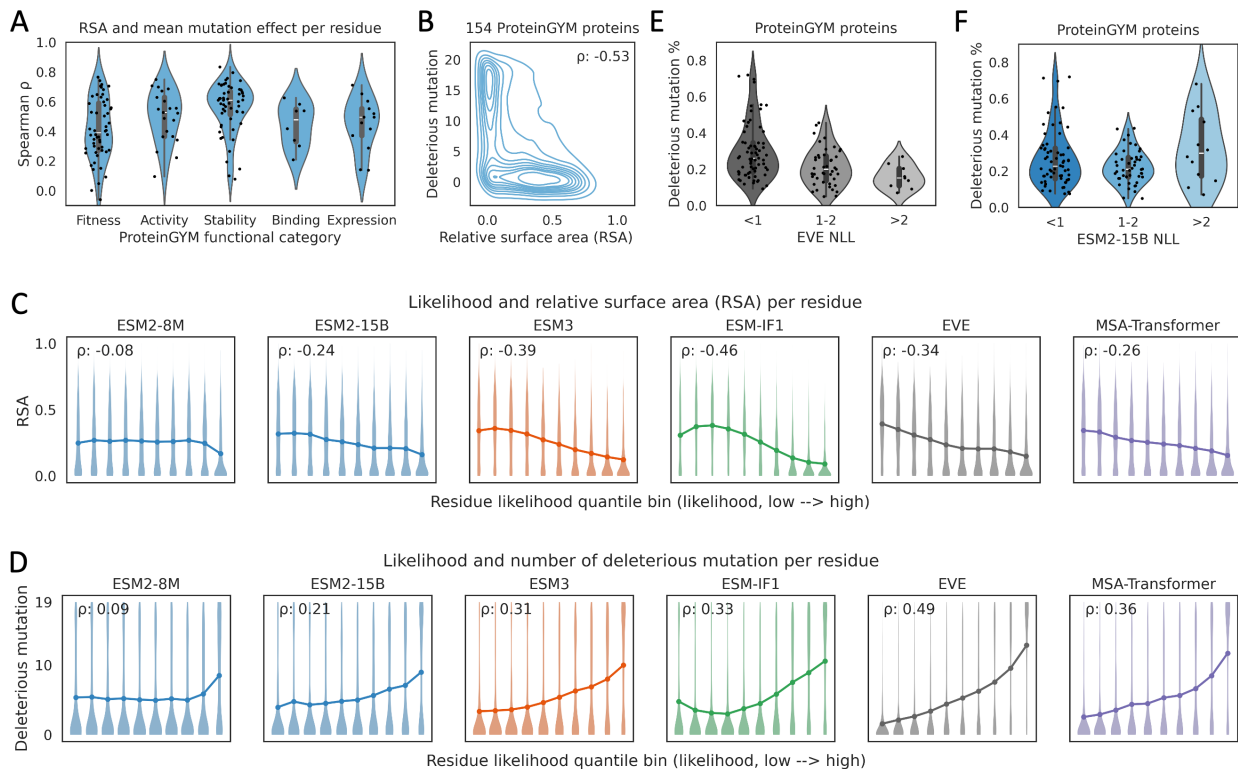
22 **Comparing model-predicted likelihood to biophysical context and mutation sensitivity.**

23 To describe the context of each residue more quantitatively, we considered both the biophysical aspect
24 and mutation sensitivity. Biophysical context is measured using relative solvent accessibility (RSA), which
25 is closely associated with mutation effects^{31,32}: mutations at buried sites tend to destabilize the protein and
26 lead to loss of function, whereas mutations at surface residues have relatively small effects on stability on
27 average, with a fraction involved in functions. We observed correlations between RSA and the mean
28 mutation effect per residue in ProteinGYM experiments, with stronger correlations observed for
29 experiments of stability (**Figure 5A**). Although the mean experimental mutation effect reflects mutation
30 sensitivity of each residue, it is not comparable across DMS experiments. Therefore, we manually
31 examined the distribution of experimental effects. Among the 154 experiments, 122 display bimodal
32 distributions (as exemplified in **Figure S9A**). For these experiments, we manually defined thresholds to
33 distinguish deleterious and neutral mutations. The thresholds were determined by visual inspection to
34 identify the separation point between two modes in the distribution of mutation effects, and the
35 corresponding threshold values are provided in the Supplementary Data 1. Mutation sensitivity at each
36 site was then quantified by counting the number of deleterious mutations among the 19 possible mutations.
37 We observed a strong inverse relationship between RSA and the number of deleterious mutations (**Figure**
38 **5B; S9B**), supporting the connection between protein structure and mutation sensitivity (functional
39 importance).

40 By comparing model-predicted per-residue likelihoods with the biophysical context and mutation
41 sensitivity, we found that residues with high predicted likelihoods by all methods tend to have lower RSA
42 and more deleterious mutations (**Figure 5C-D**). This demonstrates that both general and family-specific
43 models can capture biophysical and functional context to some extent. Biophysical context is better
44 captured by structure-informed models such as ESM3 and ESM-IF1 at high likelihoods (**Figure 5C**), owing
45 to their direct use of structural input. Mutation sensitivity, however, is better captured by family-specific
46 models (**Figure 5D, S10**). Among the top 10% of high-likelihood residues, the mean number of deleterious
47 mutations is 8.5 for ESM2-8M, 9.0 for ESM2-15B, 9.9 for ESM3, and 10.5 for ESM-IF1, compared to 12.8
48 for EVE and 11.7 for MSA-Transformer. Conversely, among the bottom 10% of low-likelihood residues,
49 the mean number of deleterious mutations is 5.3 for ESM2-8M, 3.9 for ESM2-15B, 3.3 for ESM3, and 4.9
50 for ESM-IF1, compared to 1.6 for EVE and 2.5 for MSA-Transformer (**Figure 5D**). Furthermore, at the
51 protein level, we observed that proteins with high EVE sequence likelihood tend to harbor more deleterious

1 mutations (**Figure 5E**)—a trend not observed for ESM2 (**Figure 5F**). These results indicate that general
 2 model predicted likelihood fails to reflect functional importance (mutation sensitivity) as reliably as that
 3 of family-specific models.

4



5

6 **Figure 5. Comparison of model-predicted residue likelihood with the biophysical context and mutation sensitivity.**

7 **A**, Correlations between relative solvent accessibility (RSA) and mean mutation effect per residue. Each point represents a
 8 ProteinGYM experiment. **B**, The relationship between RSA and number of deleterious mutations per residue. 122 experiments
 9 showing bimodal distributions of mutation effect are analyzed. **C-D**, Residues are grouped into ten quantile bins based on model-
 10 predicted probability. For each bin, the distribution and mean value are shown for RSA (**C**, range 0–1) and the number of deleterious
 11 mutations per residue (**D**, range 0–19). Only results for experiments with bimodal distribution of mutation effect are shown, only
 12 residues with 19 mutations are shown. **E-F**, 122 DMS experiments were classified into three groups based on sequence NLL from
 13 EVE or ESM2-15B. The proportion of deleterious mutations in each protein is shown.

14

15 **The performance of general models depends on how well predicted $p(\text{sequence})$ matches**
 16 **evolutionary patterns in homologs.**

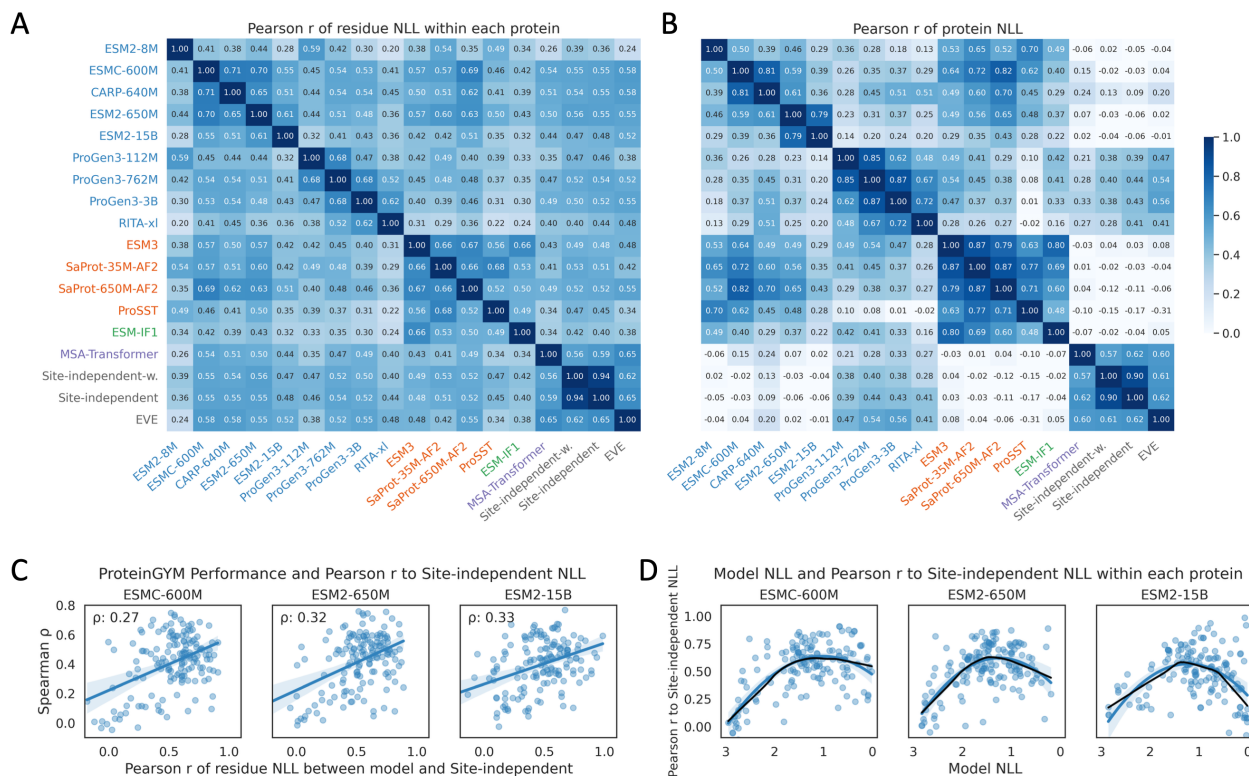
17 Finally, we set out to explain the bell-shaped relationship between sequence likelihood and fitness
 18 prediction performance. We began by comparing the predicted likelihoods from different methods. Family-
 19 specific model predicted likelihoods directly reflect evolutionary patterns in homologous sequences: a
 20 high residue-level likelihood indicates a conserved residue, while a high protein-level likelihood indicates
 21 a large fraction of conserved residues. Even though the family-specific models we evaluated employ
 22 different architectures, their predicted likelihoods are correlated at both the residue and protein levels
 23 (**Figure 6A-B**). General models, in contrast, implicitly learn evolutionary information from large-scale
 24 datasets. Likelihoods predicted by different general models are also correlated at both the residue and
 25 protein levels (**Figure 6A-B**). When comparing general models with family-specific models, correlations
 26 are observed at the residue level (**Figure 6A**), suggesting that general models capture the correct overall
 27 trend. However, for protein level likelihoods, general models show weak or no correlation with family-
 28 specific models (**Figure 6B**; decoder-based pLMs show weak protein-level correlations, see Discussion).

1 This may be because the per-residue NLL correlation is only moderate (Pearson $r \approx 0.5$; **Figure 6A**), and
 2 aggregating residue-level likelihoods to the protein level amplifies noise. This explains why ESM2-
 3 predicted likelihoods capture mutation sensitivity at the residue level (**Figure 5D**) but fail to reflect the
 4 proportion of deleterious mutations at the protein level (**Figure 5F**). Because likelihoods derived from
 5 homologs by family-specific models better reflect protein fitness (**Figure 5D-E**), we therefore hypothesized
 6 that the bell-shaped curve arises because general models capture evolutionary information most
 7 effectively when their predicted sequence likelihood falls within a moderate range.

8 To test this, we compared ESM models to family-specific models: the site-independent model, EVE, and
 9 MSA-Transformer. These models capture evolutionary patterns from sequence alone. We used likelihoods
 10 from family-specific models to represent evolutionary patterns in homologous sequences and quantified
 11 how well ESM models capture evolutionary patterns by computing the Pearson correlation between their
 12 predicted residue-level likelihoods. We found that the better ESM models capture evolutionary information,
 13 the better they perform in fitness prediction (**Figure 6C, S11A**). Notably, ESM models capture evolutionary
 14 information most strongly at moderate sequence likelihoods (**Figure 6D, S11B**), which corresponds to the
 15 range where their fitness prediction performance is optimal (**Figure 3**). These results are consistent across
 16 the three family-specific models we evaluated (**Figure 6C-D, S11A-B**). Therefore, the ability of ESM
 17 models to predict fitness relies on how well predicted likelihoods align with evolutionary patterns. These
 18 findings suggest that we should first assess whether the predicted sequence likelihoods align with
 19 evolutionary patterns, derived from family-specific models, before applying general models to fitness
 20 prediction. Considering both the fraction of proteins exhibiting this alignment and the resulting
 21 performance gains relative to no filtering, we recommend using the per-residue NLL Pearson correlation
 22 around 0.5 as the cutoff (ProteinGYM in **Figure S11C**, protein stability dataset² in **Figure S11D**).

23 We also compared other general models to family-specific models and found that ESM3 and ESM-IF1
 24 behaved differently (**Figure S12**). Unlike sequence-only pLMs, these structure-informed models can
 25 capture not only evolutionary patterns derived from homologs but also biophysical constraints (**Figure 5C**),
 26 the latter being particularly useful for modeling mutation effects on protein stability (**Table S1, Figure S13**).

27



28

1 **Figure 6. The performance of general models depends on how well they capture evolutionary information.**

2 **A**, Pearson correlation of per-residue NLL within each protein, the mean correlations across proteins in 154 experiments are shown.
3 **B**, Pearson correlation of whole-sequence NLL, which is calculated as the mean of per-residue NLL within the protein. **C**, Relationship
4 between ESM fitness prediction performance and per-residue NLL correlation to the site-independent model. The y-axes show
5 fitness prediction performance. The x-axes show the Pearson correlation between ESM NLL and site-independent NLL per residue
6 within each protein. Each point represents an experiment. The curves represent linear regressions, with the shaded areas indicating
7 the 95% confidence intervals. **D**, Relationship between ESM sequence likelihood and correlation to site-independent NLL. The y-
8 axes show the Pearson correlation between ESM NLL and site-independent NLL per residue within each protein. The x-axes show
9 the ESM2 NLL. Each point represents an experiment. The black curves show LOWESS modeling, while the colored curves represent
10 second-order polynomial regressions, with the shaded areas indicating the 95% confidence intervals.

11

12 **Discussion**

13 Predicting the protein fitness landscape is one of the most important applications of pLMs and other deep
14 learning models in biology. Our work explains why these models can predict fitness, why larger models
15 do not always perform better, and under what conditions they should or should not be used for fitness
16 prediction. Overall, we found that the better a general model captures evolutionary information, the better
17 its performance on fitness prediction. However, unlike family-specific models that directly utilize
18 evolutionary patterns in homologs, general models are influenced by unrelated factors which can decouple
19 predicted likelihoods from true evolutionary patterns. This disconnect leads to poor performance on fitness
20 prediction at extreme sequence likelihood, which leads to the bell-shaped relationship between model
21 performance and predicted likelihood. The optimal likelihood for general models corresponds to
22 $p(sequence) \approx 0.3$. Interestingly, a 30% sequence identity threshold is commonly used to search
23 homologous sequences³³. One reason for the bell-shaped relationship is the probabilistic coupling
24 between wild-type and mutant residues imposed by the SoftMax operation in these models. Because
25 supervised models usually do not use sequence likelihood to infer mutation effects and do not enforce
26 this probabilistic coupling, they are not expected to exhibit this scaling behavior. Supervised fine-tuning
27 of pLMs for fitness prediction showed that larger models tend to achieve better performance¹⁶.

28 Besides the factors we mentioned in the result section, model-predicted sequence likelihood is also
29 influenced by model architecture and training strategy. For example, models that use MSAs as input tend
30 to predict higher likelihoods, as it's easier to predict sequence with homologous sequences as input.
31 Training strategies also play a role: higher weights on the regularization loss can suppress overfitting,
32 reducing likelihoods of training sequences. Moreover, sequence patterns can also affect predicted
33 likelihoods³⁴. Therefore, predicted likelihoods should be interpreted with caution. In our evaluation, most
34 analyses were conducted within each model, so differences in model architecture and training strategy do
35 not affect our results. For comparisons between models, we used Pearson correlation which captures
36 similarity in likelihood patterns independent of absolute values.

37 We observed a monotonic relationship between model predicted residue likelihood and performance on
38 19 mutations at each residue (**Figure 4B**), we explained it as follows: residues assigned low probabilities
39 correspond to positions where general models poorly capture contextual information or where family-
40 specific models suffer from poor homolog alignment. This results in noisy predicted probabilities of the 20
41 amino acids and reduces per-residue fitness prediction performance. In contrast, residues with higher
42 predicted probabilities indicate that general models better capture contextual information, or that family-
43 specific models have good homolog alignment. As shown in Figure 5D, these residues tend to correspond
44 to functionally important positions. Figure 6A further suggests that general models identify these residues
45 by capturing the conservation signal. For these residues, models assign high probability to the wild-type
46 amino acid and low probability to alternatives, causing most mutations to be predicted as deleterious.
47 This aligns with their functional importance (mutation sensitive) and results in improved per-residue
48 performance.

49 The bell-shaped relationship is less pronounced for the decoder-based autoregressive pLMs ProGen3
50 and RITA compared with other general models (**Figure S7A**), and their peak performance at the optimal
51 likelihood remains relatively low. This likely reflects an intrinsic limitation of autoregressive pLMs for fitness
52 prediction: they cannot simultaneously incorporate contextual information from both upstream and

1 downstream residues (ESM-IF1 is also autoregressive, but it leverages complete structural context).
2 Nonetheless, autoregressive models have advantages over masked models: they can explicitly estimate
3 sequence likelihood, whereas masked language models rely on pseudo-likelihood; and they are trained
4 on all residues of training sequences, while masked models are trained only on a subset of residues
5 masked during training. These may explain why autoregressive pLMs show correlations with family-
6 specific models in protein level likelihood, whereas masked models do not (**Figure 6B**). Additionally,
7 autoregressive models are suited for generating protein sequences of varying lengths.

8 Weinstein et al.³⁵ studied the benefits of model misspecification in fitness prediction and attributed the
9 decreased performance of larger pLMs to their improved density estimation. Here, we show that
10 likelihoods predicted by pLMs do not always reliably reflect the true density of homologs (**Figure 6B**); thus,
11 the reduced performance of larger pLMs cannot be solely explained by improved density estimation.
12 Gordon et al.³⁶ reported that the preference (measured by likelihood) for a given protein sequence
13 established during pretraining is predictive of pLM fitness prediction performance. However, they did not
14 provide a mechanistic explanation. Yu et al.³⁷ found that model prediction entropy is related to viral protein
15 fitness prediction, with low entropy (corresponding to high sequence likelihood) associated with better
16 performance. Gurev et al.³⁸ reported that larger pLMs perform better on viral protein fitness prediction.
17 While their conclusions on viral proteins may appear to conflict with our findings, we note that viral proteins
18 are underrepresented in protein datasets and are therefore assigned relatively low likelihoods by general
19 models (**Figure S14A**). As a result, their sequence likelihoods have not reached the optimal levels observed
20 for well-represented proteins, thus the expected performance decline with increasing likelihood or model
21 size is not observed (**Figure S14A**). This phenomenon also extends to other proteins that are poorly
22 learned by general models, such as de novo designed proteins (**Figure S14B**).

23 Understanding model-predicted likelihood is a fundamental question in language models and the key to
24 understand how models generate their outputs. Interestingly, a bell-shaped relationship has also been
25 observed between LLM sentence likelihoods and human quality judgments³⁹. Thus, an important caveat
26 for both LLMs and protein models is that data points predicted with high likelihoods may not be real or
27 biologically meaningful. Prior work has shown that overparameterized models trained on small datasets
28 tend to memorize data: assigning high likelihoods to training data. In contrast, smaller models trained on
29 large datasets tend to generalize by learning shared patterns⁴⁰. For applying protein models to mutation
30 effect prediction, certain level of generalization, where the model can integrate information from homologs,
31 is more desirable. In our analysis, current models appear to operate in the intermediate (determined by
32 the ratio of training tokens to model parameters⁴⁰): they memorize some proteins / residues, generalize to
33 a subset, and ignore others. Which proteins / residues a model chooses to memorize, generalize, and
34 ignore remains an open question.

35 Our results offer practical guidance for applying general models to predict fitness, we recommend first
36 verifying whether the predicted likelihoods align with evolutionary patterns in homologs, which can be
37 calculated using simple family-specific models (**Figure S11C-D**). For training next generation fitness
38 predictors, we recommend incorporating evolutionary patterns during training. This can be achieved by
39 estimating evolutionary patterns from homologs beforehand and encouraging alignment between
40 predicted likelihoods and evolutionary patterns within the training objective. While our study focuses on
41 protein mutation effects, the issue we found is general and may extend to DNA/RNA language models and
42 other applications, such as protein design.

43

44 **Supplementary information**

45 Supplementary Figures S1-14

46 Supplementary Table S1

47 Supplementary Date:

- 48 1. 154 DMS experiments evaluated in this study, including manually defined cutoffs for
49 classifying deleterious and neutral mutations.

- 1 2. Performance on 154 DMS experiments across different methods.
- 2 3. Mean mutation effect prediction performance.
- 3 4. Performance per residue.
- 4 5. Predicted likelihood per residue from different methods.

5

6 **Author contributions**

7 Y.S. and C.H. conceived the study. C.H. performed the analyses and drafted the manuscript. Y.S. and
8 C.H. interpreted the results. All authors revised the manuscript.

9

10 **Competing interests**

11 The authors declare no competing interests.

12

13 **Acknowledgment**

14 This work was supported by NIH grants R35GM149527 and Simons Foundation SFARI #1019623.

15

16 **Data and code availability**

17 All mutation data used in this manuscript are publicly available from ProteinGYM (<https://proteingym.org/>)
18 and Zenodo (DOI: 10.5281/zenodo.7401274). The codes for running models are the same as that provided
19 by ProteinGYM (<https://github.com/OATML-Markslab/ProteinGym>). The site-independent models and
20 marginal approaches can be readily implemented following the method descriptions in our work.

21

22 **Methods**

23 **Dataset, protein structure, and sequence analysis.**

24 In this study, we focused exclusively on single-residue substitutions. Other mutation types—including
25 multiple-residue substitutions, insertions, deletions, and truncations—were not considered for the
26 following reasons: (1) current models' performance on these mutations is markedly lower compared to
27 single substitutions³; (2) several of the evaluated methods cannot be directly applied to these mutation
28 types; and (3) while some studies have extended models to these mutation types, the predictions are
29 derived from prediction of single-residue substitutions⁹ (e.g., the prediction of a multiple-residue
30 substitution is calculated as the sum of the corresponding single-substitution predictions, and the
31 prediction of a truncation is calculated as the maximum single-substitution prediction in the truncated
32 region).

33 Experimental mutation effects, model predictions, Multiple sequence alignments (MSAs), and protein
34 structures were downloaded from the ProteinGYM³ website (<https://proteingym.org/>) and its GitHub
35 repository (<https://github.com/OATML-Markslab/ProteinGym>) in May 2025. From the 217 mutational
36 scanning datasets available, only residues with experimental effects of all 19 possible substitutions were
37 included. Restricting to proteins with at least 20 such residues resulted in a dataset comprising 486,932
38 single-residue substitution mutations across 154 experiments (Supplementary Data 1). We note that this
39 ProteinGYM filtering criterion was applied to all analyses except for the PTEN analysis, where we included
40 residues with at least ten measured substitutions to increase residue coverage. The thresholds for
41 classifying deleterious and neutral mutations provided in ProteinGYM are not used, as they are defined
42 based on the distribution of mutation effects of both single and multiple substitutions and appear
43 unreasonable for some proteins in our visualizations. Instead, we manually defined thresholds using only

1 single substitutions. The cutoffs were determined by visual inspection to identify the separation point
2 between the two modes in the distribution of experimental single-mutation effects. Cutoffs were defined
3 only for DMS experiments exhibiting a clear bimodal separation, as exemplified in Figure S9A. The
4 corresponding cutoff values are provided in Supplementary Data 1 to ensure reproducibility.

5 The mega-scale protein folding stability dataset², including both designed and natural mini domains, was
6 analyzed in Figure S4E-H, S6, S11D, and S14B. Single-residue substitutions in the table
7 “Tsuboyama2023_Dataset2_Dataset3_20230416.csv” were analyzed. For applying the site-independent
8 model to these proteins in Figure S11D, homologs were identified using the ColabFold⁴¹ server with default
9 parameters, only homologs with less than 50% gaps were analyzed.

10 PTEN cancer hotspot mutations were downloaded from cBioPortal²⁵ (February 2026). The curated set of
11 non-redundant studies was used, missense mutations annotated as “CancerHotspot” or “3DHotspot”
12 were included.

13 ProteinGYM provides predictions from a wide range of methods, only likelihood-based methods were
14 evaluated here. For models (EVE, MSA-Transformer, ESM1v) offering both single-model and ensemble
15 predictions, the single-model predictions were used. Protein solvent-accessible surface areas (SASA)
16 were calculated using the MDTraj⁴² package (version: 1.10.2). Relative surface areas were computed by
17 dividing the SASA of each residue by the maximum SASA observed for that amino acid across all proteins.
18 For homolog detection, MMseqs2⁴³ (version 16.747c6) was used to search against UniRef50²⁰ (version
19 2021_04, which was used to train ESM2) with the following parameters: `--min-seq-id 0.2 -c 0.8 --max-
20 accept 1000`.

21

22 **Applying masked language models for protein fitness prediction.**

23 Masked language models (MLMs) evaluated in the study include ESM2¹⁰, ESM1v²², ESMC²⁷, CARP²⁸,
24 ESM3-open¹³ (1.4B parameters), SaProt²¹, ProSST³⁰ (110M parameters, 2048 structure tokens), and MSA-
25 Transformer¹¹. Among them, ESM2, ESM1v, ESMC, and CARP are sequence-only MLMs; ESM3, ProSST,
26 and SaProt are sequence and structure hybrid MLMs; MSA-Transformer is an MSA-based MLM. For MLM,
27 computing the exact sequence likelihood $p(\mathbf{x})$ is intractable. Instead, pseudo-likelihood is used, where
28 each residue is masked individually and predicted in turn. For a protein of length L with sequence $\mathbf{x} =$
29 (x_1, \dots, x_L) (x_i denotes the wild-type amino acid, which is also explicitly written as $x_{i,wt}$ when it needs to be
30 distinguished from mutations in the following sections.), structure \mathbf{s} (structure token of the protein is also
31 a list of tokens with length L), and a set of homologs \mathbf{h} , the pseudo-likelihood of sequence-only MLM is
32 defined as:

$$33 \quad \log(\hat{p}_{sequence\ MLM}(\mathbf{x})) = \frac{1}{L} \sum_{i=1}^L \log(p(x_i | \mathbf{x}_{\setminus i})) \quad (1)$$

34 Where $\mathbf{x}_{\setminus i}$ denotes the sequence with position i masked.

35 ESM3, ProSST, and SaProt were trained with both sequence and structure tokens. ESM3 predicts the
36 probabilities of sequence and structure tokens separately, ProSST only predicts masked sequence during
37 training, whereas SaProt predicts them jointly. Specifically, SaProt predicts the probabilities of 400
38 combinations of sequence–structure tokens, corresponding to 20 amino acids \times 20 FoldSeek⁴⁴ structure
39 tokens. During fitness inference, all structure tokens were provided, while the residue at the mutation site
40 was masked. For SaProt, the predicted probability for each amino acid was obtained by summing over 20
41 sequence–structure combinations with that amino acid. The pseudo-likelihood of sequence and structure
42 hybrid MLM is defined as:

$$43 \quad \log(\hat{p}_{hybrid\ MLM}(\mathbf{x})) = \frac{1}{L} \sum_{i=1}^L \log(p(x_i | \mathbf{x}_{\setminus i}, \mathbf{s})) \quad (2)$$

1 MSA-Transformer leverages homologous sequences to provide evolutionary context. The choice of
 2 homologs can be made by random sampling, selecting the most similar, or selecting the most dissimilar
 3 sequences¹¹. In this study, we followed the protocol used in ProteinGYM, where homologs were sampled
 4 according to sequence-similarity-based weights (also used for EVE). The pseudo-likelihood of MSA-
 5 Transformer is defined as:

$$6 \quad \log(\hat{p}_{\text{MSA-Transformer}}(\mathbf{x})) = \frac{1}{L} \sum_{i=1}^L \log(p(x_i | \mathbf{x}_{\setminus i}, \mathbf{h})) \quad (3)$$

7 For fitness prediction, the log-likelihood ratio (LLR) between the mutated sequence \mathbf{x}_{mt} and the wild-type
 8 sequence \mathbf{x}_{wt} can be computed using the pseudo-likelihood of their full sequences:

$$9 \quad LLR_{\text{pseudo-likelihood}} = \log(\hat{p}(\mathbf{x}_{mt})) - \log(\hat{p}(\mathbf{x}_{wt})) \quad (4)$$

10 However, this approach is computationally intensive, since for each mutated sequence of length L ,
 11 computing the pseudo-likelihood requires L forward passes through the model, making it impractical for
 12 large-scale fitness inference. Instead, simplified approximations are used, including the masked marginal
 13 and wild-type marginal approaches, both of which consider only the predicted probabilities at the mutation
 14 site rather than the likelihood of the entire sequence. The underlying hypothesis is that the difference in
 15 predicted probabilities at a single residue is proportional to the difference in overall sequence likelihood.

16 For the masked marginal approach, the residue at the mutation site is masked, and the model directly
 17 predicts the probabilities of the wild-type and mutant amino acids, the LLR is then calculated as:

$$18 \quad LLR_{\text{masked-marginal}} = \log(p(x_{i,mt} | \mathbf{x}_{\setminus i})) - \log(p(x_{i,wt} | \mathbf{x}_{\setminus i})) \quad (5)$$

19 Previous studies have shown that for MLMs trained with the standard masking scheme (e.g., ESM2, where
 20 15% of tokens are randomly selected during pre-training, with 80% replaced by the mask token, 10%
 21 kept unchanged, and 10% replaced by a random token), predictions from the unmasked wild-type
 22 sequence can be used directly for fitness inference⁹. In this case, the LLR is calculated as:

$$23 \quad LLR_{\text{wild-type-marginal}} = \log(p(x_{i,mt} | \mathbf{x})) - \log(p(x_{i,wt} | \mathbf{x})) \quad (6)$$

24 The computational cost differs substantially among these methods. For a protein of length L , scoring all
 25 possible single-residue substitutions (19 per position) requires running the model $L \times 19 \times L$ times with
 26 $LLR_{\text{pseudo-likelihood}}$, L times with $LLR_{\text{masked-marginal}}$, and only one time with $LLR_{\text{wild-type-marginal}}$. In this
 27 study, we used masked marginal approach for all MLM methods, which is the most used methods to
 28 estimate fitness from MLM. For MLMs trained with a fraction of tokens substituted and a fraction of tokens
 29 unchanged, the sequence likelihood computed without masking is highly correlated with that from the
 30 masked approach¹⁹ (Figure S7B), and the LLR scores derived from both approaches show nearly identical
 31 performance in predicting mutation effects⁹ (Figure S7C). Therefore, our conclusions also apply when
 32 using MLMs with the wild-type marginal approach (Figure S7D).

33 For all models, the predicted probability of all possible tokens in each position sum to one:

$$34 \quad p(x_{i,wt}) + \sum_{a \neq wt} p(x_{i,a}) = 1 \quad (7)$$

35 Where a represent amino acids different from the wild-type. ProteinGYM provides LLRs for the MLMs we
 36 evaluated, but not the underlying sequence likelihoods. However, given the definition of LLRs (Eq. 5)
 37 together with Eq. 7, the predicted residue probabilities can be inferred using the following equations:

$$38 \quad \frac{1}{p(x_{i,wt})} = 1 + \sum_{a \neq wt} \frac{p(x_{i,a})}{p(x_{i,wt})} = 1 + \sum_{a \neq wt} e^{LLR_a} \quad (8)$$

$$\log(p(x_{i,wt})) = -\log\left(1 + \sum_{a \neq wt} e^{LLR_a}\right) \quad (9)$$

Eq. 8 and 9 were used to estimate sequence likelihoods for positions with LLRs of all 19 possible mutations in the ProteinGYM dataset, eliminating the need to run these models. As ProteinGYM and the default ProSST implementation use the wild-type marginal approach, masked marginal scores were computed for ProSST to ensure comparability of sequence likelihoods across methods. Sequence likelihoods of five ESM1v models and ESM2 prediction for the mega-scale protein folding stability dataset² were computed by us.

Applying autoregressive models for protein fitness prediction.

We evaluated several autoregressive generative models: pLMs RITA-xl²⁹ (1.2B parameters) and ProGen3¹⁶, and the inverse-folding model ESM-IF1¹². Both RITA and ProGen3 are decoder-only transformer models, with ProGen3 incorporating mixture-of-experts layers. These models were trained in an autoregressive manner: RITA and ProGen3 on both forward and reverse sequences, and ESM-IF1 on forward sequences only (N-terminal to C-terminal). For autoregressive models, the likelihood of the forward sequence can be computed explicitly as:

$$\log(p(\mathbf{x})) = \frac{1}{L} \sum_{i=1}^L \log(p(x_i | \mathbf{x}_{<i})) \quad (10)$$

Where $\mathbf{x}_{<i}$ denotes the residues before position i . For autoregressive pLMs, the likelihood of a sequence was computed as the mean of predictions from the forward and reverse sequences:

$$\log(p_{\text{autoregressive pLM}}(\mathbf{x})) = \frac{1}{2L} \sum_{i=1}^L [\log(p(x_i | \mathbf{x}_{<i})) + \log(p(x_i | \mathbf{x}_{>i}))] \quad (11)$$

For ESM-IF1, the likelihood of a sequence was computed as:

$$\log(p_{\text{ESM-IF1}}(\mathbf{x})) = \frac{1}{L} \sum_{i=1}^L \log(p(x_i | \mathbf{x}_{<i}, \mathbf{s})) \quad (12)$$

For fitness prediction using autoregressive models, the sequence likelihood is computed directly for each mutated sequence. Since the wild-type sequence likelihood is constant, calculating LLR is unnecessary for performance evaluation. ProteinGYM provides negative log likelihood (NLL) for each mutated sequence rather than LLRs. The sequence likelihoods for wild-type proteins were computed by us using Eq. 11 and 12.

We also investigated whether autoregressive pLMs could be applied using marginal approaches, which can reduce the times of model running from 2× the number of mutant sequences to only two (forward and reverse wild-type sequences). In this setting, the probabilities of both the wild-type and mutant amino acids at a given position are predicted using upstream and/or downstream residues as context. Using ProGen3-3B, we evaluated four strategies: (1) computing the LLR from the forward sequence only, (2) computing the LLR from the reverse sequence only, (3) averaging the LLRs obtained from the forward and reverse sequences, and (4) averaging the predicted probabilities from the forward and reverse sequences before computing the LLR. The results are summarized in Table S1. Performances of all marginal approaches are substantially lower than using the ProGen3-3B full sequence likelihood (**Table S1**). Thus, in the Results section, only the results obtained using the full sequence likelihood for generative pLMs were reported. We note that the marginal approaches for autoregressive pLMs also exhibit weak bell-shaped relationships between wild-type sequence likelihood and fitness prediction performance (data not shown).

1

2 **Site-independent models.**

3 The site-independent model was used as the baseline family-specific model. MSAs provided by
4 ProteinGYM were used, and homologous sequences with more than 50% gaps were removed. For each
5 site, the probabilities of the 20 amino acids and the gap character were calculated as:

$$6 \quad p(x_i = a) = \frac{\sum_h w_h \delta(x_i^h, a) + \epsilon}{\sum_h w_h + 21\epsilon} \quad (13)$$

7 where w_h denotes the weight of homologous sequence h , which was either calculated using the same
8 strategy as in EVE or set as one for all sequences. δ is the Kronecker delta. ϵ is the pseudo count, for the
9 unweighted model, ϵ was set to 1; for the weighted model, ϵ was set to the minimum positive value in the
10 frequency table. Fitness was estimated using LLR of mutant and wild-type amino acids at each position.
11 Benchmarking on the ProteinGYM mutations evaluated in this study, both the site-independent
12 (unweighted) and the site-independent-weight model achieve mean Spearman correlations around 0.36
13 (**Table S1**). Despite their simplicity, they outperform the site-independent model trained with EVmutation⁴⁵
14 (mean Spearman correlation 0.33).

15

16 **Family-specific model EVE and its marginal approach.**

17 EVE⁸ models were trained separately for each MSA using the code provided in the ProteinGYM GitHub
18 repository, taking approximately a week on ten A6000 GPUs. Model was trained with the random seed of
19 0 and `threshold_focus_cols_frac_gaps = 1` (all positions were modeled). As a variational autoencoder⁴⁶,
20 EVE does not provide exact sequence likelihoods. Instead, the evidence lower bound (ELBO) was used
21 as a proxy, calculated as:

$$22 \quad ELBO = \mathbb{E}_{q(z|x)}[\log(p(x|z))] - KL[q(z|x) || p(z)] \quad (14)$$

23 Where z is the latent variable, $q(z|x)$ is the approximate posterior from the encoder, $p(x|z)$ is the
24 decoder likelihood, and $KL[\cdot||\cdot]$ denotes the Kullback–Leibler divergence (KLD). The ELBO was computed
25 for both wild-type and mutated sequences by averaging over 20,000 samples of z . The decoder
26 conditional probability $p(x|z)$ was used to approximate the probability of each amino acid at each site,
27 averaged over the 20,000 samples using wild-type sequence as input.

28 Using EVE to estimate fitness is computationally intensive, as it requires training a separate model with
29 millions of parameters for each MSA, and during inference, the decoder must be run tens of thousands of
30 times per mutated sequence to obtain a stable ELBO estimation. To reduce this burden, we implemented
31 a simplified approach, which we termed EVE-marginal. In this approach, the conditional probability $p(x|z)$
32 is obtained from the wild-type sequence only, and fitness is estimated like the MLM wt-marginal
33 approach. EVE-marginal performs comparable to the full EVE model (**Table S1**), while being substantially
34 efficient.

35

36 **Reference**

- 37 1. Fowler, D. M. & Fields, S. Deep mutational scanning: a new style of protein science. *Nat. Methods* **11**, 801–807
38 (2014).
39 2. Tsuboyama, K. *et al.* Mega-scale experimental analysis of protein folding stability in biology and design. *Nature*
40 **620**, 434–444 (2023).
41 3. Notin, P. *et al.* ProteinGym: Large-Scale Benchmarks for Protein Fitness Prediction and Design. *Adv. Neural Inf.*
42 *Process. Syst.* **36**, 64331–64379 (2023).
43 4. Cheng, J. *et al.* Accurate proteome-wide missense variant effect prediction with AlphaMissense. *Science*
44 <https://doi.org/10.1126/science.adg7492> (2023) doi:10.1126/science.adg7492.

- 1 5. Zhang, H., Xu, M. S., Fan, X., Chung, W. K. & Shen, Y. Predicting functional effect of missense variants using graph
2 attention neural networks. *Nat Mach Intell* **4**, 1017–1028 (2022).
- 3 6. Gao, H. *et al.* The landscape of tolerated genetic variation in humans and primates. *Science* **380**, eabn8153 (2023).
- 4 7. Livesey, B. J. & Marsh, J. A. Using deep mutational scanning to benchmark variant effect predictors and identify
5 disease mutations. *Mol. Syst. Biol.* **16**, e9380 (2020).
- 6 8. Frazer, J. *et al.* Disease variant prediction with deep generative models of evolutionary data. *Nature* **599**, 91–95
7 (2021).
- 8 9. Brandes, N., Goldman, G., Wang, C. H., Ye, C. J. & Ntranos, V. Genome-wide prediction of disease variant effects
9 with a deep protein language model. *Nat Genet* **55**, 1512–1522 (2023).
- 10 10. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science* **379**,
11 1123–1130 (2023).
- 12 11. Rao, R. *et al.* MSA Transformer. <https://doi.org/10.1101/2021.02.12.430858> (2021)
13 doi:10.1101/2021.02.12.430858.
- 14 12. Hsu, C. *et al.* Learning inverse folding from millions of predicted structures. *bioRxiv* 2022.04.10.487779 (2022)
15 doi:10.1101/2022.04.10.487779.
- 16 13. Hayes, T. *et al.* Simulating 500 million years of evolution with a language model. *Science* **387**, 850–858 (2025).
- 17 14. Xiong, J. *et al.* Guide your favorite protein sequence generative model. Preprint at
18 <https://doi.org/10.48550/arXiv.2505.04823> (2025).
- 19 15. Kaplan, J. *et al.* Scaling Laws for Neural Language Models. Preprint at <https://doi.org/10.48550/arXiv.2001.08361>
20 (2020).
- 21 16. Bhatnagar, A. *et al.* Scaling unlocks broader generation and deeper functional understanding of proteins. *bioRxiv*
22 2025.04.15.649055 (2025) doi:10.1101/2025.04.15.649055.
- 23 17. Notin, P. Have We Hit the Scaling Wall for Protein Language Models? *Pascal Notin*
24 <https://pascalnotin.substack.com/p/have-we-hit-the-scaling-wall-for> (2025).
- 25 18. Chen, B. *et al.* xTrimoPGLM: unified 100-billion-parameter pretrained transformer for deciphering the language of
26 proteins. *Nat. Methods* **22**, 1028–1039 (2025).
- 27 19. Gordon, C., Lu, A. X. & Abbeel, P. Protein Language Model Fitness Is a Matter of Preference. *bioRxiv*
28 2024.10.03.616542 (2024) doi:10.1101/2024.10.03.616542.
- 29 20. UniProt, C. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res* **47**, D506–D515 (2019).
- 30 21. Su, J. *et al.* Democratizing protein language model training, sharing and collaboration. *Nat. Biotechnol.*
31 <https://doi.org/10.1038/s41587-025-02859-7> (2025) doi:10.1038/s41587-025-02859-7.
- 32 22. Meier, J. *et al.* Language models enable zero-shot prediction of the effects of mutations on protein function. *bioRxiv*
33 2021.07.09.450648 (2021) doi:10.1101/2021.07.09.450648.
- 34 23. Mighell, T. L., Evans-Dutson, S. & O’Roak, B. J. A Saturation Mutagenesis Approach to Understanding PTEN Lipid
35 Phosphatase Activity and Genotype-Phenotype Relationships. *Am. J. Hum. Genet.* **102**, 943–955 (2018).
- 36 24. Matreyek, K. A., Stephany, J. J., Ahler, E. & Fowler, D. M. Integrating thousands of PTEN variant activity and
37 abundance measurements reveals variant subgroups and new dominant negatives in cancers. *Genome Med.* **13**,
38 165 (2021).
- 39 25. Cerami, E. *et al.* The cBio cancer genomics portal: an open platform for exploring multidimensional cancer
40 genomics data. *Cancer Discov.* **2**, 401–404 (2012).
- 41 26. Cleveland, W. S. Robust Locally Weighted Regression and Smoothing Scatterplots. *J. Am. Stat. Assoc.* **74**, 829–
42 836 (1979).
- 43 27. Evolutionary Scale · ESM Cambrian: Revealing the mysteries of proteins with unsupervised learning.
44 <https://www.evolutionaryscale.ai/blog/esm-cambrian>.
- 45 28. Yang, K. K., Fusi, N. & Lu, A. X. Convolutions are competitive with transformers for protein sequence pretraining.
46 *Cell Syst.* **15**, 286–294.e2 (2024).
- 47 29. Hesslow, D., Zanichelli, N., Notin, P., Poli, I. & Marks, D. RITA: a Study on Scaling Up Generative Protein Sequence
48 Models. Preprint at <https://doi.org/10.48550/arXiv.2205.05789> (2022).
- 49 30. Li, M. *et al.* ProSST: Protein Language Modeling with Quantized Structure and Disentangled Attention. Preprint at
50 <https://doi.org/10.1101/2024.04.15.589672> (2024).
- 51 31. Zhong, G., Zhao, Y., Zhuang, D., Chung, W. K. & Shen, Y. PreMode predicts mode-of-action of missense variants
52 by deep graph representation learning of protein sequence and structural context. *Nat. Commun.* **16**, 7189 (2025).
- 53 32. Tsishyn, M., Hermans, P., Rooman, M. & Pucci, F. Residue conservation and solvent accessibility are (almost) all
54 you need for predicting mutational effects in proteins. *Bioinformatics* **41**, btaf322 (2025).
- 55 33. Pearson, W. R. An introduction to sequence similarity (“homology”) searching. *Curr Protoc Bioinforma.* **Chapter 3**,
56 3.1.1–3.1.8 (2013).
- 57 34. Kantroo, P., Wagner, G. P. & Machta, B. B. In-Context Learning can distort the relationship between sequence
58 likelihoods and biological fitness. Preprint at <https://doi.org/10.48550/arXiv.2504.17068> (2025).

- 1 35. Weinstein, E. N., Amin, A. N., Frazer, J. & Marks, D. S. Non-identifiability and the Blessings of Misspecification in
2 Models of Molecular Fitness. 2022.01.29.478324 Preprint at <https://doi.org/10.1101/2022.01.29.478324> (2023).
- 3 36. Gordon, C. W., Lu, A. X. & Abbeel, P. Protein Language Model Fitness is a Matter of Preference. in (2024).
- 4 37. Yu, Y., Jiang, F., Zhong, B., Hong, L. & Li, M. Entropy-driven zero-shot deep learning model selection for viral
5 proteins. *Phys. Rev. Res.* **7**, 013229 (2025).
- 6 38. Gurev, S., Youssef, N., Jain, N. & Marks, D. S. Variant effect prediction with reliability estimation across priority
7 viruses. 2025.08.04.668549 Preprint at <https://doi.org/10.1101/2025.08.04.668549> (2025).
- 8 39. Zhang, H., Duckworth, D., Ippolito, D. & Neelakantan, A. Trading Off Diversity and Quality in Natural Language
9 Generation. Preprint at <https://doi.org/10.48550/arXiv.2004.10450> (2020).
- 10 40. Morris, J. X. *et al.* How much do language models memorize? Preprint at
11 <https://doi.org/10.48550/arXiv.2505.24832> (2025).
- 12 41. Mirdita, M. *et al.* ColabFold: making protein folding accessible to all. *Nat Methods* **19**, 679–682 (2022).
- 13 42. McGibbon, R. T. *et al.* MDTraj: A Modern Open Library for the Analysis of Molecular Dynamics Trajectories. *Biophys*
14 *J* **109**, 1528–32 (2015).
- 15 43. Steinegger, M. & Soding, J. MMseqs2 enables sensitive protein sequence searching for the analysis of massive
16 data sets. *Nat Biotechnol* **35**, 1026–1028 (2017).
- 17 44. van Kempen, M. *et al.* Fast and accurate protein structure search with Foldseek. *Nat. Biotechnol.*
18 <https://doi.org/10.1038/s41587-023-01773-0> (2023) doi:10.1038/s41587-023-01773-0.
- 19 45. Hopf, T. A. *et al.* Mutation effects predicted from sequence co-variation. *Nat Biotechnol* **35**, 128–135 (2017).
- 20 46. Kingma, D. P. & Welling, M. Auto-Encoding Variational Bayes. Preprint at <https://doi.org/10.48550/arXiv.1312.6114>
21 (2022).
- 22