PROTMAMBA: A HOMOLOGY-AWARE BUT ALIGNMENT-FREE PROTEIN STATE SPACE MODEL

Anonymous authors

Paper under double-blind review

Abstract

Protein design has important implications for drug discovery, personalized medicine, and biotechnology. Models based on multiple sequence alignments efficiently capture the evolutionary information in homologous protein sequences, but multiple sequence alignment construction is imperfect. We present ProtMamba, a homology-aware but alignment-free protein language model based on the Mamba architecture. In contrast with attention-based models, ProtMamba efficiently handles very long context, comprising hundreds of protein sequences. We train ProtMamba on a large dataset of concatenated homologous sequences, using two GPUs. We combine autoregressive modeling and masked language modeling through a fill-in-the-middle training objective. This makes the model adapted to various protein design applications. We demonstrate ProtMamba's usefulness for the generation of novel sequences and for fitness prediction. ProtMamba reaches competitive performance with other protein language models despite its smaller size, which sheds light on the importance of long-context conditioning.

027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 Proteins are essential building blocks of life, serving vital roles in metabolic processes, cellular 029 transport, structural integrity, and immune responses. Composed of long chains of amino acids (polypeptides), proteins fold into specific three-dimensional structures critical for their biological 031 functions. One of the key challenges in biology is protein engineering and design: conceiving protein sequences to exhibit enhanced or novel functions. While experimental approaches like directed 033 evolution and mutational scanning are effective in this regard, they only allow exploring the neighbors 034 of existing sequences. However, the recent growth of extensive databases has opened up new avenues for computational methods that exploit the breadth of biological evolution. For instance, UniProt (The UniProt Consortium, 2021) contains more than two hundreds of millions of protein sequences. 036 Biological functions exert evolutionary constraints on protein sequences, which can be probed by 037 considering families of homologous proteins (i.e. proteins that share an evolutionary history) and analyzing this data through statistical methods and, more recently, through deep learning methods.

Protein language models rely on recurrent (Bepler & Berger, 2019), transformer (Rives et al., 2021) 040 or convolutional (Yang et al., 2024) architectures, and are trained through masked language modeling, 041 autoregressive modeling, or discrete diffusion techniques (Alamdari et al., 2023), on large ensembles 042 of single protein sequences (Khakzad et al., 2023). The representations learned by these models 043 correlate with biochemical properties of proteins (such as function, structure, contacts) (Elnaggar 044 et al., 2021; Vig et al., 2021; Rives et al., 2021; Madani et al., 2023), and can be used to generate 045 protein sequences or to evaluate the fitness of variants. The vast majority of these methods are 046 trained on non-structured ensembles of single protein sequence and do not have direct access to 047 homology, or to conservation and variability within protein families. Models trained on multiple 048 sequence alignments (MSAs) of homologous sequences have also been introduced, despite raising memory challenges and potentially suffering from the imperfections of MSAs (Thompson et al., 2011). Successful MSA-based transformer models, such as MSA Transformer (Rao et al., 2021) 051 or the EvoFormer module of AlphaFold2 (Jumper et al., 2021) alternate attention along protein sequences and across homologs. More recently, PoET (Truong Jr & Bepler, 2024) was trained on 052 concatenations of non-aligned homologous sequences, offering a promising autoregressive alternative to MSA Transformer for protein fitness prediction and design.

State space models such as S4 (Gu et al., 2021), Hyena (Poli et al., 2023) and Mamba (Gu & Dao, 2023) are catching up with transformers thanks to their ability to efficiently handle very long sequences of tokens. These models were quickly adapted to work with biological data. Approaches such as HyenaDNA (Nguyen et al., 2023) or Evo (Nguyen et al., 2024) were trained on long DNA sequences and capture regulatory mechanics. Meanwhile, PTM-Mamba addresses post-translational modifications of protein sequences (Peng et al., 2024).

060 In this paper, we present ProtMamba, a novel homology-aware but alignment-free protein language 061 model, trained on concatenated sequences of homologous proteins. Based on the Mamba architecture 062 (Gu & Dao, 2023), ProtMamba is able to handle extremely long contexts (unlimited lengths during 063 inference). Trained to autoregressively predict the next amino acid, but also with a fill-in-the-middle 064 (FIM) objective, it can be used for multiple different tasks. First, ProtMamba can autoregressively generate novel sequences without contextual information. Second, by providing ProtMamba with 065 sequences from a specific protein family or subfamily as context, users can prompt it to generate 066 sequences tailored to their specifications. This conditional generation approach is a key strength 067 of the model (see also Truong Jr & Bepler (2024)), and could become an alternative to fine-tuning. 068 Third, ProtMamba supports sequence inpainting, i.e., filling specific masked regions with the desired 069 number of amino acids. For this, along with homologous sequences (used as context), the model is provided with a target sequence to be modified. This generation mode opens novel methods of 071 designing specific parts of protein sequences. Furthermore, ProtMamba is useful for fitness prediction tasks. Users can input a sequence with specific masked positions, prompting the model to output 073 the probability distribution of all mutations in each variant with a single forward pass. Across these 074 various tasks, we obtain competitive results with larger protein language models and task-specific 075 methods.

076 077

078 079 080

081 082

084

085

090

092

093

095

096

2 Methods

2.1 Key technical contributions

1. To harness the evolutionary information present in homologous sequences without relying on multiple sequence alignments (MSAs), we use as input a concatenation of homologous sequences for each protein family. In each of these long arrays, sequences are separated with a specific token. The motivation is that evolutionary information is extremely useful for protein modeling Jumper et al. (2021); Rao et al. (2021); Abramson et al. (2024), but MSAs can be inaccurate. This approach is similar to that used recently in the autoregressive transformer PoET (Truong Jr & Bepler, 2024).

2. We develop an architecture based on Mamba blocks, an alternative to attention recently proposed by (Gu & Dao, 2023) that relies on state space models. In Mamba, which is a recurrent neural network, time complexity scales linearly in sequence length, bypassing the quadratic time complexity constraints of transformers. This allows handling significantly longer input sequences, in addition to being faster to train and to use at inference. This is a key asset here, as concatenating homologous sequences results in long inputs. Note that Truong Jr & Bepler (2024) employed attention matrix chunking to address this issue, but this results in potential losses of statistical dependence signals, and only partially solves the memory limit.

- 3. We combine elements of both autoregressive modeling and masked language modeling 097 (MLM), by training our model using the fill-in-the-middle (FIM) objective (Bavarian et al., 098 2022; Fried et al., 2022; Raffel et al., 2020). The model learns to predict masked patches 099 extracted randomly from a sequence and positioned at the end of it, and can therefore leverage the full sequence context, while being trained autoregressively. This is of particular interest for biological sequences, because preceding and subsequent tokens can all be informative to 102 predict a new token. While autoregressive models are generative by definition, they yield 103 the probability of each new token conditioned on previous ones (ignoring subsequent ones). Besides, MLM can be productively used for protein sequence generation (Sgarbossa et al., 105 2023).
- 4. To promote the model's ability to reason over in-sequence positions, which is particularly useful for the FIM task, we modify the original Mamba implementation by introducing sequence-level positional embeddings. This enables the model to pay attention to relative

positions inside each sequence. In inference and generation, it opens the possibility of controlling the number of amino-acids to generate.

110 111

108

2.2 MODEL ARCHITECTURE AND TRAINING STRATEGY

ProtMamba's architecture is adapted from Mamba (Gu & Dao, 2023). An important modification is that we introduce learned positional embeddings for the input tokens. Among different variants (see supplementary section C), we observed that the most effective and stable method to integrate positional embeddings is to concatenate them with the input token embeddings into a single vector. Specifically, we allocated half of the embedding dimension *d* to token information and the other half to positional information.

119 We trained a 107 million parameters model with 16 layers, embedding dimension d = 1024, and 120 hidden state dimension equal to embedding dimension. We started with a maximal total input sequence 121 length of $2^{11} = 2048$ amino acids (recall that input sequences are concatenated homologous protein 122 sequences). The model was trained following (Gu & Dao, 2023) with some minor modifications. 123 We used the AdamW optimizer with the following parameter values: weight decay w = 0.1 and 124 $(\beta_1, \beta_2) = (0.9, 0.95)$. We scheduled the learning rate to increase from zero to 6×10^{-4} with a linear warm-up of 500 steps followed by a constant learning rate. To optimize memory usage, we trained 125 the model using the bfloat16 format. 126

To avoid training instabilities observed in (Nguyen et al., 2024; Waleffe et al., 2024), we implemented
a callback mechanism to revert to a previous checkpoint if the loss never assumed values below a
threshold for 10 successive evaluation steps. The threshold value was chosen as the lowest training
loss increased by 0.5%. This ensures that the loss decreases overall, while allowing it to transiently
increase. We also prevented gradient explosions by clipping the gradient norm to 1.0.

132 The model was trained by scheduling the context length of the input using sequence length warm up 133 (SLW) (Nguyen et al., 2023). Initially, we used inputs of length $L = 2^{11}$ tokens with a batch size of 134 64. We doubled input length each time the loss reached a plateau, simultaneously reducing batch size 135 to maintain a fixed total number of tokens per batch. In case of memory constraints, we decrease 136 the batch size and use gradient accumulation. This heuristic approach is based on the idea that a 137 longer context should provide more information. It is useful because of training instabilities for long contexts (Nguyen et al., 2023; 2024). Note that we did not start training the model with a long context 138 to benefit from a larger batch size, which helps to approximate the loss landscape more efficiently. 139 Finally, once we reached a context length of $L = 2^{17}$, we implemented gradient checkpointing to 140 minimize memory consumption. This allowed us to increase the batch size for the final part of the 141 training and obtain a better approximation of the loss landscape, see (Nguyen et al., 2023; 2024). 142

The model was trained on one NVIDIA RTX A6000 GPU for 35 days, and then on two of them for 143 15 days. This allowed us to keep the batch size large enough when the context size increased. In total, 144 the model was trained on 1.95×10^{11} tokens (approximately 1.5 epochs) and used 2.0×10^{20} FLOPs 145 during training. These numbers show the huge improvements that the Mamba architecture has in 146 terms of training speed with respect to transformers. As a comparison, the smallest ESM3 model 147 (Hayes et al., 2024) was trained with 0.8×10^{11} tokens using 6.72×10^{20} FLOPs, which means that 148 given a fixed amount of compute, ProtMamba can see 8.5 times the tokens seen by ESM3. See Figure 149 S1 for the training curves. 150

We consider two different ProtMamba versions that were obtained by saving checkpoints at different moments of the training. Our model *ProtMamba, Foundation* was trained on a maximum context length of 2¹⁵ tokens. Our model *ProtMamba Long, Foundation* was trained until the context length reached 2¹⁷ tokens. Both models were fine-tuned for 2 days on predicting only the FIM amino acids to improve inpainting capabilities, yielding the models *ProtMamba/ProtMamba Long, Fine-tuned*.

- We also performed multiple ablations on the model architecture and on the different modalities by training models on 10B tokens for 50k steps (see supplementary section C).
- 158
- 159 2.3 DATASET CONSTRUCTION
- 161 We trained ProtMamba on OpenProteinSet (Ahdritz et al., 2024), a dataset which comprises 16 millions MSAs, one for each sequence cluster within Uniclust30 (Mirdita et al., 2017). This dataset

162 was curated to train OpenFold (Ahdritz et al., 2022). We used a filtered subset of the full dataset, con-163 sisting of maximally diverse representative MSA clusters, built by iteratively eliminating redundant 164 clusters whose representative sequences appeared in other clusters' MSAs (Ahdritz et al., 2024). This 165 ensures that each representative sequence is only present in its cluster, as detailed in (Ahdritz et al., 166 2024). This dataset comprises 268,000 clusters including a total of 508 million sequences and 110 billion residues (see Figure S2 for additional statistics). A validation set and a testing set are formed 167 by holding out respectively 192 and 500 randomly chosen clusters from the training set. Importantly, 168 our use of the filtered version of OpenProteinSet (Ahdritz et al., 2024) ensures that overlap between clusters in the training, validation and test set is strongly minimized. Indeed, this filtering is based on 170 selecting only MSAs of maximal diversity and ensuring that the reference sequences used to build 171 each cluster are not present in any other cluster. 172

Figure 1 illustrates the construction of a training example. First, a cluster is randomly selected from 173 the filtered OpenProteinSet database described in Section 2.3. As OpenProteinSet uses MSAs, we 174 restore the original unaligned sequences by removing gaps and converting all lowercase insertion 175 residues to uppercase. Each amino acid is tokenized using a unique token. Then, N sequences are 176 sampled uniformly at random and concatenated into a single array, with a **<cls>** token separating 177 each sequence from the next one. The value of N is chosen for the total length of the concatenated 178 sequence to exceed the desired training context length L (e.g. $L = 2^{11}$ at the beginning of training), 179 and the input is then cropped precisely at L. Next, the sequences are prepared for the FIM task. For 180 each sequence, some patches of consecutive tokens are randomly sampled (see below) and masked 181 by replacing them with a mask token **<mask i>**, with one such token representing patch *i*. For each 182 patch, we append to the sequence another mask token followed by the corresponding masked amino 183 acids (which are unmasked). An <eos> token is used to separate the main (masked) sequence from its unmasked patches.



199 Figure 1: Input to ProtMamba. Each element of the input is a concatenation of unaligned homolo-200 gous sequences separated by **<cls>** tokens. Each sequence starts with a **<cls>** token and ends with an 201 <eos> token. Masked segments are replaced by numbered mask tokens, <mask1>, ..., <mask5>. 202 The masked tokens are appended to the sequence, after the **<eos>** token, each masked segment being 203 preceded by its associated mask token. The position indices ("pos-ids") follow the succession of 204 tokens in the natural sequence. Thus, the masked tokens have their initial position indices in the 205 natural sequence. The position index of each mask is set to that of the first associated masked token. In this particular example we sampled two masks i = 1, 2 with length $P_1 = 3$ and $P_2 = 2$. 206

207 208

209

210

211

212

213

The following rules are applied when masking each sequence:

- The number of masked patches in a sequence is sampled from a Poisson distribution with λ = 1, and capped at 5 (by resampling in case values above 5 are obtained). This yields no mask in 36% of sequences, one mask in 36% of sequences, and more in 28% of sequences.
 The starting position of each patch is sampled uniformly (without replacement) from all
- possible positions in the sequence.
- 3. The length P_i of each patch *i* is sampled uniformly in $[1, \max(P_i)]$, where $\max(P_i)$ is 0.2 times the distance from the start point of patch *i* to the start point of patch *i* + 1 (or to the end of the sequence for the last patch). This ensures that no more than 20% of all tokens

in each sequence are masked, in line with masking fractions of similar models (Rao et al., 2021; Rives et al., 2021).

Finally, each token is allocated a position index (used to obtain the associated positional embedding) that tracks its position in the original sequence. The position indices of **<cls>** and **<eos>** are set to zero, while the mask tokens **<mask i>** have the same position indices as the first token they are masking, see Figure 1.

3 Results

216

217

218

224 225

226 227

228

229

230

231

232

233

234

235

236

237

238

3.1 PROTMAMBA BENEFITS FROM LONG CONTEXT

To evaluate the effectiveness of incorporating context information in ProtMamba, we examine the scaling of the model's perplexity with context length for natural sequences. Perplexity is commonly used to evaluate autoregressive models and assesses how uncertain they are about a sequence. It is the exponential of the cross entropy loss. Figure 2 shows the scaling of perplexity for the masked parts of the sequences as a function of the number of context sequences, when using the FIM objective. ProtMamba Long (Fine-tuned) achieves remarkably low values of perplexity for small numbers N_m of masked tokens. Furthermore, perplexity decreases when increasing the number of context sequences, revealing the positive impact of richer context on model performance. This decrease tends to be steeper for larger N_m , suggesting that these difficult tasks particularly benefit from richer context. Given the diverse lengths of sequences across protein families, we report perplexity versus the number of sequences in the context rather than versus the total length of the context. Indeed, there can be different amounts of information in contexts of similar lengths but composed of sequences of varying lengths.



Figure 2: Scaling of the FIM perplexity with the number of context sequences. We show the FIM perplexity for different numbers N_m of masked amino acids versus the number of context sequences. Results are averaged over all 500 clusters of the test set and 100 replicates for each cluster (differing by the random sampling of context sequences). Context sizes go up to 2^{17} amino acids. To reduce noise, we take the exponential moving average, and we restrict to cases where the count of samples is at least 100. See Figure S3 for a log-log version of this figure.

260 Furthermore, we study the scaling of the per-sequence perplexity (i.e. the standard autoregressive 261 perplexity of the full non-masked sequence) computed on the test set using ProtMamba Long 262 (Foundation), see Figure S4. Initially, we notice a decrease of perplexity to a minimum of 7.70 as 263 the number of sequences in the context increases, with lower perplexity values for shorter individual 264 sequences, but this reduction plateaus after a certain point. We attribute this behavior to the finite 265 size (d = 1024), see Section 2) of the hidden state of the model, which limits its capacity to 266 effectively leverage context information at each step. We hypothesize that a larger model with 267 a higher-dimensional hidden state could increase the amount of information transferred from the context to the next predicted token. For completeness, we also report perplexity versus context length 268 measured in tokens (see Figure S5). There, we observe a rise in perplexity when the context sizes 269 reaches $2^{17} = 131,072$ tokens, which is the highest context length seen during training. We expect that further training the model for longer contexts could lead to lower perplexity values, yet ultimately
 reaching a lower bound due to the limitations imposed by the hidden state dimension and model size.

273 274

3.2 PROTMAMBA PREDICTS MUTATIONAL EFFECTS IN DIFFERENT PROTEIN FAMILIES

275 Next, we evaluate ProtMamba's ability to predict mutational effects, leveraging its inpainting ca-276 pabilities arising from the FIM training objective. Indeed, by masking specific amino acids in the 277 wild-type sequence of interest, we can predict the fitness of all variants at these sites. Our first step 278 to evaluate variant fitness is to collect a context of homologs to the wild-type sequence. We use the ColabFold protocol (Mirdita et al., 2022) for this, ensuring that diverse sequences are found in a few 279 minutes. Then, we randomly subsample 200 sequences among those that have between 30% and 280 98% similarity to the wild type to construct the context, and we sort these sequences by increasing 281 similarity to the wild type, as in Truong Jr & Bepler (2024). 282

To evaluate the effect of a variant with a single mutated site, we append the wild-type sequence to the context, mask the mutated residue in it, and predict this residue using the FIM method. Let C denote the union of the context sequences and of the wild-type sequence masked at the mutated position *i*. We evaluate the effect of mutations at position *i* by their fitness score \mathcal{F} , defined as:

 $\mathcal{F}(i, x_i, \mathcal{C}) = \log p(x_i | \mathcal{C}) - \log p(x_i^{WT} | \mathcal{C}),$

288

289 for all residues x_i different from the wild-type residue x_i^{WT} . Using this method based on the FIM 290 objective allows us to evaluate the effects of all mutations at position i, decreasing 20-fold the number 291 of passes through the model needed to evaluate all mutations with respect to the typical method used with autoregressive protein language models. To predict the fitness effects of variants involving 292 293 mutations at multiple sites, we add all the single mutation likelihoods. This approximate, but fast method avoids computing the complete likelihood for all variants, thus reducing the number of calls 294 to ProtMamba. It is accurate when the mutations can be considered independent. We also test variant 295 scoring by ProtMamba using the autoregressive log-likelihood loss instead of the FIM loss (this 296 approach is called "ProtMamba AR"). Details on the different approaches we employed to predict 297 mutational effects with ProtMamba are given in Supplementary section B. 298

299 We consider the ProteinGym benchmark (Notin et al., 2023), which contains 217 datasets of substitutions in protein sequences (both single and multiple) and allows comparing to state-of-the-art 300 methods. In Table 1, we report the performance of ProtMamba, and we compare it to published 301 models classified by type: alignment-based models, single sequence protein language models (PLMs), 302 aligment-enhanced PLMs, homology-aware PLMs, and structure-aware models. Table 1 shows 303 that variant scoring by ProtMamba using FIM outperforms using the autoregressive log-likelihood 304 (ProtMamba AR). All ProtMamba performances reported in Table 1 are those of ProtMamba Long 305 fine-tuned on the FIM task, except for ProtMamba AR where we used ProtMamba Long. Indeed, 306 considering the 4 different ProtMamba versions (see Section 2), we found that ProtMamba models 307 fine-tuned on the FIM task outperform foundation models, and that ProtMamba Long performs better 308 than ProtMamba, confirming the importance of training the model with a long context (see Figure S6). 309 In Figure S7(a), we break down the performance of ProtMamba Long for different context lengths and different protein sequences lengths. We observe that variants with long sequences particularly 310 benefit from long contexts, as they allow including more sequences. This interpretation is supported 311 by Figure S7(b), which shows that this dependence on context length is weaker when considering 312 context length in terms of number of sequences. Based on performance on a validation set (see 313 supplementary Section A and Figure S8), we chose to use a context of 200 sequences to predict 314 fitness using ProtMamba Long (fine-tuned). 315

Table 1 shows that ProtMamba outperforms single-sequence PLMs ("PLM" type) of the same size
 (ESM-2, 150M), and performs similarly or better than larger models like Tranception L and ESM-2
 (650M). This illustrates the power of homology information for mutational effect prediction.

Since MSA information remains useful in scoring variants, in the rows "Alignment + PLM" of Table
1, we show results where explicit use of MSAs was made, either via retrieval, i.e. ensembling the
models with an independent-site model, as in Notin et al. (2023) (denoted by "R"), or by combining
a PLM with GEMME in Marquet et al. (2024). Using retrieval, ProtMamba (w/ R) obtains similar
performance as Tranception L (w/ R) and as MSA Transformer, which leverage MSA information.
These two models were trained using more than one order of magnitude more FLOPs than ProtMamba.

Model type	Model	#params	ρ	Time	Citation
Alignment-based	Site-Independent	-	0.359	-	Hopf et al. (2017)
	GEMME	-	0.455	-	Laine et al. (2019)
PLM	Tranception L (w/o R)	700M	0.374	-	Notin et al. (2022)
	ESM-2	150M	0.387	-	Lin et al. (2023)
	ESM-2	650M	0.414	-	Lin et al. (2023)
Homology-aware	ProtMamba (single)	107M	0.406	7m	
PLM	ProtMamba AR (single)	107M	0.367	1h 39m	
	PoET (single)	201M	0.447	9h 51m	Truong Jr & Bepler (2024)
	PoET (ensemble)	201M	0.470	148h*	Truong Jr & Bepler (2024)
Alignment	ProtMamba (w/ R)	107M	0.432	10m	
+ PLM	MSA-Transformer	100M	0.421	-	Rao et al. (2021)
	Tranception L (w/ R)	700M	0.434	-	Notin et al. (2022)
	VespaG	3B	0.458	-	Marquet et al. (2024)
Structure-aware	ESM-IF1	142M	0.422	-	Hsu et al. (2022)
	SaProt	650M	0.457	-	Su et al. (2023)
	ProSST	110M	0.507	-	Li et al. (2024)

Table 1: Performance of ProtMamba and of existing models on the ProteinGym benchmark. For 340 different models classified by type, and whose numbers of parameters are given, we show Spearman correlation ρ values between predicted and experimentally measured variant effects in ProteinGym. 342 We denote the MSA-augmented methods with retrieval by "(w/R)". New models introduced here are 343 highlighted in bold. They include ProtMamba with and without retrieval, where variants are scored 344 using the FIM loss, and ProtMamba AR, where they are scored using the autoregressive log-likelihood. 345 Published results were obtained from https://proteingym.org/. For ProtMamba and PoET, 346 we also report the time needed to score all variants in ProteinGym (excluding homolog retrieval). Top 347 performances in terms of ρ and time are highlighted in gray. 348

Estimated as $15 \times$ the time taken by PoET (single).

349 350

339

341

In Table 1, the state-of-the-art model on the ProteinGym benchmark is ProSST (Li et al., 2024), 351 a structure-aware model. Among structure-agnostic models, the state-of-the-art model is PoET, a 352 homology-aware transformer. ProtMamba reaches a performance which is only slightly lower than 353 PoET. This is notable, as PoET has twice more parameters, and is much slower at scoring variants 354 than ProtMamba. Table 1 shows that ProtMamba can score all ProteinGym variants in \sim 7 to 10 355 minutes on a single NVIDIA RTX A6000 GPU, while PoET takes ~ 10 hours for this, and up to ~ 6 356 days in the ensemble mode (the top structure-agnostic method), using the same hardware. 357

In Table S1, we break down results by MSA depth and by number of mutations. In Figure S9, we 358 further break down the comparisons between models on ProteinGym by category of experiment 359 (panel (a)), taxonomic category (panel (b)) and sequence length (panel (c)). We also show scores for 360 different models on randomly selected example experimental datasets in Figure S10. 361

362 363

364

3.3 PROTMAMBA ACCURATELY PREDICTS THE ACTIVITY OF CHORISMATE MUTASE ENZYMES

Next, we evaluate ProtMamba, and in particular the power of the FIM objective, on a dataset of 365 experimentally tested natural and *in silico* generated sequences from the chorismate mutase family 366 from Russ et al. (2020). Chorismate mutase functions as an enzyme involved in the catalysis 367 of synthesis of amino acids, and is a domain of the bifunctional chorismate mutase/prephenate 368 dehydratase. We use ProtMamba to evaluate the activity of experimentally studied variants of this 369 enzyme. For this, we sample 100 sequences, either randomly among all natural sequences that were 370 experimentally studied, or randomly among the subset of those that were experimentally shown to be 371 active in *Escherichia coli*. For these two types of context, we test three different protocols to predict 372 the activity of the other variants in the dataset of Russ et al. (2020) with ProtMamba. First, we use 373 only the chorismate mutase domains (cropped sequences) as context, and autoregressively evaluate the 374 likelihood of the full sequence ("from left to right"). Second, we use the full sequences (chorismate 375 mutase/prephenate dehydratase) as context and we evaluate the perplexity of the full sequence autoregressively from left to right. Third, we use the full sequences (chorismate mutase/prephenate 376 dehydratase) as context and evaluate the perplexity of the chorismate mutase domain using the FIM 377 objective.

378		Domain	Drotoin	Drotoin FIM	# seq.	# residues	ρ
379	Dublished methods	Domain	Flotem	Floteni, Flivi	0	0k	0.31
380	DCA energy	0.41	_	_	5	1k	0.46
381	Logistic Regression	0.43	-	-	10	3k	0.48
382	ProtMamba				20 50	7k	0.51
383	Context: any variant	0.41	0.44	0.46	50 100	19K 38k	0.49
384	Context: active variants	0.50	0.52	0.53	200	77k	0.52

387

388

389

390

Table 2: Activity prediction of chorismate mutase variants. (Left) For published methods (Russ et al., 2020), and for ProtMamba with various context types (rows) and protocols (columns, see main text), we report the Spearman correlation ρ between experimental activity and predictions. Associated ROC curves are shown in Figure S11. (Right) Effect of increasing context size (number of sequences and corresponding total number of residues) on Spearman correlation ρ for ProtMamba predictions, using only active variants in context, full domains and FIM.

391 392

The left panel of Table 2 provides a comparison of ProtMamba against published methods (Russ et al., 2020) for the two different context types and the three different protocols. We observe that using only active variants in the context consistently improves the predictive power of ProtMamba. With both context types, using full sequences is better than using only domains, and using FIM improves accuracy. Note in addition that using FIM reduces the computation time per variant compared to autoregressively scoring the full sequence.

399 The right panel of Table 2 shows the impact of context size on the performance of our best ProtMamba-400 based activity predictor (using only active variants in context, full domains and FIM). The accuracy 401 of this predictor initially strongly increases with context length, and then appears to plateau from 402 a context length around 75 sequences or 28k residues (see also Figure S12 a). Thus, context size 403 plays a critical role for these activity predictions. In Figure S12 b and c, we further compare the 404 perplexity of the variants, when using only active variants as context, and when using both inactive 405 and active variants as context. We observe that the perplexity of inactive variants is often higher when 406 using a context of active variants, indicating a better ability to predict inactivity. Furthermore, the 407 perplexity of active variants is often lower in this case, showing a better ability to predict activity with high-quality context. Additionally, we display the distribution of the perplexity for ProtMamba using 408 FIM and a context composed of active variants, compared to the experimental activity in Figure S13, 409 and the same perplexity versus the model score from Russ et al. (2020) in Figure S14. 410

- 411
- 412 413

3.4 PROTMAMBA AUTOREGRESSIVELY GENERATES PROMISING NOVEL SEQUENCES

Finally, we evaluate ProtMamba on the autoregressive generation of novel protein sequences given a 414 context of known homologs, corresponding to members of a given cluster of sequences. We generate 415 sequences from 19 randomly selected clusters in the test set, varying the following parameters: 416 temperature (T), top-k number, and top-p fraction, following the approach proposed by (Ferruz et al., 417 2022). These parameters are commonly employed to control the output of autoregressive models. At 418 each step, top-k limits their output to the top-k most probable tokens, while top-p only includes the 419 top tokens reaching a cumulative probability p. Meanwhile, temperature T adjusts the randomness 420 of sampling. Additionally, we vary the number of sequences in the context to assess the impact 421 of different levels of conditioning on the generated sequences. Specifically, for each cluster, we perform generation using context lengths of n = 10, 100, 500, 1000 and N sequences, where 422 N is the total number of sequences in the cluster. For each value of n, we consider the following 423 (T, top-k, top-p) triplets: (0.8, 10, 0.9), (0.9, 10, 0.95), (1, 10, 0.95), (1, 10, 1), (1, 15, 1). We gen-424 erate 100 sequences for each (n, T, top-k, top-p), obtaining a total of 2500 sequences per family. 425 As expected, we observed that the parameters which promote higher sampling variability tend to 426 yield sequences with higher perplexity. Note that sequences with more than 750 amino acids, i.e. 427 longer than the longest natural sequence considered here, were discarded from further analysis. They 428 represented $\sim 5\%$ of the generated sequences. 429

We compare the sampled sequences (aggregated across all parameter sets mentioned above) with
 natural sequences from the cluster used as context for generation using various scores evaluating
 novelty, homology, and structure.



Figure 3: **Comparison of low-perplexity generated sequences with natural ones.** We report the median and the standard deviation of sequence length, Hamming distance to the closest natural neighbor in the sequence cluster from which the context is drawn ("Min Hamming"), HMMER score (rescaled), pLDDT and pTM scores from ESMFold. For each of 19 test clusters, we compare the 100 sequences with lowest perplexity values out of 2500 generated sequences (*x*-axis) with a randomly chosen subset of 100 natural sequences in the sequence cluster (*y*-axis). Dashed black lines: y = x.

- 1. Novelty is assessed by computing the pairwise Hamming distance (using pairwise Smith-Waterman alignment) with each natural sequence in the cluster, after which it is possible to focus on distance to the closest natural neighbor if desired.
- 2. Homology evaluation involves training an HMM (using HMMER (Eddy, 2020)) on the cluster's MSA, obtained from OpenProteinSet, and computing the scores it gives to generated sequences.
- 3. Structure is assessed by predicting the structure of each sampled sequence using ESMFold (Lin et al., 2023). As ESMFold is a single sequence model, it provides predictions that are less biased by MSAs than those of MSA-based models. Futhermore, it is faster than AlphaFold2. ESMFold's confidence measures, both global with pTM scores and local with pLDDT scores, allow for a precise comparison of different sequences sampled from the same cluster.

We observe that ProtMamba's estimated sequence perplexity correlates well with HMMER scores, Hamming distance to the closest natural neighbor in the cluster and structural scores (see Figure S15). Thus, ProtMamba assigns lower perplexity values to sequences that are more likely to be part of the cluster. The absolute Pearson correlation value averaged over all clusters and scores is above 0.57. Detailed results for each family and each score are presented in Figure S16. Figure 3 shows that the median scores of our generated sequences that have low perplexity are comparable to those of natural ones. Overall, these results are promising for protein design applications.

Model	ProtMamba	EvoDiff-MSA	MSA Trans.	Potts	Natural
pLDDT (\uparrow) scPerplexity (\downarrow)	$egin{array}{c c} 0.75 \pm 0.13 \\ 2.63 \pm 0.45 \end{array}$	$0.60 \pm 0.16 \\ 3.17 \pm 0.58$	$0.54 \pm 0.18 \\ 3.37 \pm 0.64$	$0.56 \pm 0.14 \\ 3.17 \pm 0.51$	$\begin{vmatrix} 0.77 \pm 0.13 \\ 2.66 \pm 0.49 \end{vmatrix}$

Table 3: **Performance of ProtMamba and other models at homolog-conditioned generation.** We report two structural scores, namely the pLDDT from ESMFold (Lin et al., 2023) and the scPerplexity from ProteinMPNN (Dauparas et al., 2022) for a set of 250 protein sequences generated using ProtMamba, each from a different cluster in our test set. Note that scPerplexity is the self-consistency Perplexity computed by ProteinMPNN from the ESMFold structures obtained for each generated sequence. We compare these values to those obtained for 250 protein sequences generated by EvoDiff-MSA, MSA-Transformer and Potts models, retrieved from the Zenodo archive associated to the EvoDiff paper (Alamdari et al., 2023), and which were generated each from a different cluster of the EvoDiff validation set. We also compare to a subset of the same size of natural sequences sampled from the same test set clusters as for ProtMamba.

 \uparrow (resp. \downarrow) indicates that higher (resp. lower) scores are better.

Finally, in Table 3, we compare the generative ability of ProtMamba to that of other models that can
 perform sequence generation conditioned on homologs from a specific protein family. Following
 the approach of Alamdari et al. (2023), we randomly sample 250 clusters from our test set, and, for

486 each cluster, we generate a sequence using ProtMamba conditioned on homologs randomly sampled 487 from the cluster. We then compare two structural scores of these generated sequences with those 488 obtained for sequences generated by EvoDiff-MSA (Alamdari et al., 2023), MSA Transformer (Rao 489 et al., 2021) and Potts models (Russ et al., 2020) and provided in Alamdari et al. (2023), and for 490 natural sequences. We find that ProtMamba outperforms existing models on the homolog-conditioned generation task, and generates sequences that obtain scores comparable to those of natural sequences. 491 Note that the Hamming distance to the closest natural homolog of the sequences sampled using 492 ProtMamba is 0.56 ± 0.10 , similar to the value obtained for natural sequences (0.48 ± 0.17) and not 493 smaller, consistently with Figure 3. 494

4 DISCUSSION

496 497 498

495

Here, we presented ProtMamba, a homology-aware but alignment-free generative protein language
 model. ProtMamba leverages the long-context capabilities of state space models, allowing it to handle
 concatenated sequences of homologous proteins. It also benefits from their faster speed compared to
 attention-based models (Gu & Dao, 2023), allowing fast sequence generation and mutational effect
 prediction. ProtMamba was trained using a hybrid strategy combining autoregressive modeling and
 masked language modeling via the FIM objective. This allows ProtMamba to efficiently predict the
 next amino acid in a protein sequence as well as to inpaint masked regions.

505 Our results demonstrate ProtMamba's versatility across multiple tasks, including conditioned genera-506 tion and protein fitness prediction, both for close and for distant variants. For homolog-conditioned 507 generation, ProtMamba outperforms the state-of-the-art model EvoDiff-MSA (Alamdari et al., 2023). 508 For fitness prediction, the sequence inpainting abilities of ProtMamba, via the FIM objective, proved 509 to be particularly useful. Indeed, this functionality allows the model to exploit the full sequence context, without restricting to previous tokens as with autoregressive generation. This allows Prot-510 Mamba to reach similar performance levels as larger models, in a fraction of the time. Overall, 511 ProtMamba benefits from capturing signal across multiple scales. In particular, it is able to predict 512 fitness by exploiting constraints shared broadly across the proteome via its pre-training, but also 513 specific constraints shared between homologs via the context, and it can exploit the full context of a 514 given protein sequence when predicting only part of it. 515

516

Limitations. So far, ProtMamba did not reach perplexity values as low as those of larger transformer
 models like PoET (Truong Jr & Bepler, 2024) for full sequences. However, it can handle longer
 context sizes and requires much shorter training and inference times, which is extremely beneficial
 for the sequence inpainting task. We believe that scaling the model to larger sizes and training times
 (comparable to PoET) may result in comparable performance, while retaining ProtMamba's assets of
 lower memory cost and inference time.

We did not provide a direct test of the generative ability of ProtMamba for protein sequence inpainting.
Indeed, this is a highly specific task lacking clear benchmarks so far. However, we believe that our two analyses on fitness prediction constitute a convincing indirect proof of the usefulness of ProtMamba's inpainting ability. It would be very interesting to experimentally test ProtMamba's inpainting ability, as well as its de novo sequence generation ability (Verkuil et al., 2022).

528

Perspectives. Our results demonstrate ProtMamba's flexibility, as it allows for precise conditioning by carefully choosing the context information (e.g. restricting to active sequences). Thus, ProtMamba responds very well to prompt engineering. We propose that this could become an alternative or complement to fine-tuning of language models. ProtMamba is also naturally designed to take advantage of retrieval augmented generation (RAG) techniques (Lewis et al., 2021), as it allows for using retrieved protein sequences from any external database, to condition the generation process.

Furthermore, we envision the possibility to use the model for homology search, by scoring sequences within specific contexts. This would be very fast, because only one forward pass would be required.

An interesting further extension of ProtMamba would be to make it explicitly structure-aware, e.g.
using a structural alphabet (van Kempen et al., 2023), along the lines of SaProt (Su et al., 2023) or
ProstT5 (Heinzinger et al., 2023). Another possible extension would be to include Gene Ontology
(GO) terms to condition sequence generation (Madani et al., 2023; Nijkamp et al., 2023).

540 5 **REPRODUCIBILITY STATEMENT** 541

We describe in details all the steps to reproduce our work in Section 2 and we provide all the code in the supplementary material attached to the submission.

References

542

543

544

546

571 572

573

577

578

579

- 547 J. Abramson, J. Adler, J. Dunger, R. Evans, T. Green, A. Pritzel, O. Ronneberger, L. Willmore, 548 A. J. Ballard, J. Bambrick, S. W. Bodenstein, D. A. Evans, C. C. Hung, M. O'Neill, D. Reiman, K. Tunyasuvunakool, Z. Wu, A. e, E. Arvaniti, C. Beattie, O. Bertolli, A. Bridgland, A. Cherepanov, 549 M. Congreve, A. I. Cowen-Rivers, A. Cowie, M. Figurnov, F. B. Fuchs, H. Gladman, R. Jain, Y. A. 550 Khan, C. M. R. Low, K. Perlin, A. Potapenko, P. Savy, S. Singh, A. Stecula, A. Thillaisundaram, 551 C. Tong, S. Yakneen, E. D. Zhong, M. Zielinski, A. dek, V. Bapst, P. Kohli, M. Jaderberg, 552 D. Hassabis, and J. M. Jumper. Accurate structure prediction of biomolecular interactions with 553 AlphaFold 3. Nature, May 2024. 554
- Gustaf Ahdritz, Nazim Bouatta, Sachin Kadyan, Qinghui Xia, William Gerecke, Timothy J O'Donnell, 555 Daniel Berenberg, Ian Fisk, Niccolò Zanichelli, Bo Zhang, et al. OpenFold: Retraining AlphaFold2 556 yields new insights into its learning mechanisms and capacity for generalization. Biorxiv, pp. 2022–11, 2022. 558
- 559 Gustaf Ahdritz, Nazim Bouatta, Sachin Kadyan, Lukas Jarosch, Dan Berenberg, Ian Fisk, Andrew Watkins, Stephen Ra, Richard Bonneau, and Mohammed AlQuraishi. OpenProteinSet: Training 561 data for structural biology at scale. Advances in Neural Information Processing Systems, 36, 2024.
- 562 Sarah Alamdari, Nitya Thakkar, Rianne van den Berg, Alex X. Lu, Nicolo Fusi, Ava P. Amini, and 563 Kevin K. Yang. Protein generation with evolutionary diffusion: sequence is all you need. *bioRxiv*, 564 2023. doi: 10.1101/2023.09.11.556673. URL https://www.biorxiv.org/content/ 565 early/2023/09/12/2023.09.11.556673. 566
- Mohammad Bavarian, Heewoo Jun, Nikolas Tezak, John Schulman, Christine McLeavey, Jerry 567 Tworek, and Mark Chen. Efficient training of language models to fill in the middle. arXiv, 2022. 568
- 569 Tristan Bepler and Bonnie Berger. Learning protein sequence embeddings using information from 570 structure. In International Conference on Learning Representations, 2019. URL https:// openreview.net/forum?id=SygLehCqtm.
- J. Dauparas, I. Anishchenko, N. Bennett, H. Bai, R. J. Ragotte, L. F. Milles, B. I. M. Wicky, A. Courbet, R. J. de Haas, N. Bethel, P. J. Y. Leung, T. F. Huddy, S. Pellock, D. Tischer, F. Chan, 574 B. Koepnick, H. Nguyen, A. Kang, B. Sankaran, A. K. Bera, N. P. King, and D. Baker. Robust 575 deep learning-based protein sequence design using proteinmpnn. Science, 378(6615):49-56, 2022. 576 doi: 10.1126/science.add2187. URL https://www.science.org/doi/abs/10.1126/ science.add2187.
 - Sean R. Eddy. HMMER: biosequence analysis using profile hidden Markov models, 2020. URL http://hmmer.org.
- Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Yu Wang, Llion Jones, Tom 581 Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, Debsindhu Bhowmik, and Burkhard 582 Rost. ProtTrans: Towards cracking the language of life's code through self-supervised deep 583 learning and high performance computing. IEEE Transactions on Pattern Analysis and Machine 584 Intelligence, pp. 1-1, 2021. doi: 10.1109/TPAMI.2021.3095381. 585
- 586 Noelia Ferruz, Steffen Schmidt, and Birte Höcker. ProtGPT2 is a deep unsupervised language model for protein design. *Nature communications*, 13(1):4348, 2022.
- 588 Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, 589 Wen-tau Yih, Luke Zettlemoyer, and Mike Lewis. Incoder: A generative model for code infilling 590 and synthesis. arXiv preprint arXiv:2204.05999, 2022.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. arXiv, 592 (arXiv:2312.00752), 2023. doi: 10.48550/arXiv.2312.00752. URL http://arxiv.org/abs/ 2312.00752.

- Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- Thomas Hayes, Roshan Rao, Halil Akin, Nicholas J. Sofroniew, Deniz Oktay, Zeming Lin, Robert Verkuil, Vincent Q. Tran, Jonathan Deaton, Marius Wiggert, Rohil Badkundri, Irhum Shafkat, Jun Gong, Alexander Derry, Raul S. Molina, Neil Thomas, Yousuf Khan, Chetan Mishra, Carolyn Kim, Liam J. Bartie, Matthew Nemeth, Patrick D. Hsu, Tom Sercu, Salvatore Candido, and Alexander Rives. Simulating 500 million years of evolution with a language model. *bioRxiv*, 2024. doi: 10.1101/2024.07.01.600583. URL https://www.biorxiv.org/content/early/2024/07/02/2024.07.01.600583.
- Michael Heinzinger, Konstantin Weissenow, Joaquin Gomez Sanchez, Adrian Henkel, Milot Mirdita,
 Martin Steinegger, and Burkhard Rost. Bilingual language model for protein sequence and structure.
 bioRxiv, 2023. doi: 10.1101/2023.07.23.550085.
- Thomas A Hopf, John B Ingraham, Frank J Poelwijk, Charlotta PI Schärfe, Michael Springer,
 Chris Sander, and Debora S Marks. Mutation effects predicted from sequence co-variation. *Nat Biotechnol*, 35(2):128, 2017.
- Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and Alexander Rives. Learning inverse folding from millions of predicted structures. *bioRxiv*, 2022. doi: 10.1101/2022.04.10.487779. URL https://www.biorxiv.org/content/early/ 2022/09/06/2022.04.10.487779.
- J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool,
 R. Bates, A. Žídek, A. Potapenko, A. Bridgland, C. Meyer, S. A. A. Kohl, A. J. Ballard, A. Cowie,
 B. Romera-Paredes, S. Nikolov, R. Jain, J. Adler, T. Back, S. Petersen, D. Reiman, E. Clancy,
 M. Zielinski, M. Steinegger, M. Pacholska, T. Berghammer, S. Bodenstein, D. Silver, O. Vinyals,
 A. W. Senior, K. Kavukcuoglu, P. Kohli, and D. Hassabis. Highly accurate protein structure
 prediction with AlphaFold. *Nature*, 596(7873):583–589, 2021. doi: 10.1038/s41586-021-03819-2.
- Hamed Khakzad, Ilia Igashov, Arne Schneuing, Casper Goverde, Michael Bronstein, and Bruno
 Correia. A new age in protein design empowered by deep learning. *Cell Systems*, 14(11):
 925–939, 2023. ISSN 2405-4712. doi: 10.1016/j.cels.2023.10.006. URL https://www.
 sciencedirect.com/science/article/pii/S2405471223002983.
- Elodie Laine, Yasaman Karami, and Alessandra Carbone. Gemme: A simple and fast global epistatic model predicting mutational effects. *Molecular Biology and Evolution*, 36(11):2604–2619, 08 2019. ISSN 0737-4038. doi: 10.1093/molbev/msz179. URL https://doi.org/10.1093/molbev/msz179.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks, 2021.
- Mingchen Li, Pan Tan, Xinzhu Ma, Bozitao Zhong, Huiqun Yu, Ziyi Zhou, Wanli Ouyang, Bingxin
 Zhou, Liang Hong, and Yang Tan. Prosst: Protein language modeling with quantized structure
 and disentangled attention. *bioRxiv*, 2024. doi: 10.1101/2024.04.15.589672. URL https:
 //www.biorxiv.org/content/early/2024/05/17/2024.04.15.589672.1.
- ⁶³⁷ Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023. doi: 10.1126/ science.ade2574.
- A. Madani, B. Krause, E. R. Greene, S. Subramanian, B. P. Mohr, J. M. Holton, J. L. Olmos, C. Xiong, Z. Z. Sun, R. Socher, J. S. Fraser, and N. Naik. Large language models generate functional protein sequences across diverse families. *Nat. Biotechnol.*, 41(8):1099–1106, 2023.
- 646 Celine Marquet, Julius Schlensok, Marina Abakarova, Burkhard Rost, and Elodie Laine. Vespag:
 647 Expert-guided protein language models enable accurate and blazingly fast fitness prediction. *bioRxiv*, pp. 2024–04, 2024.

648 Milot Mirdita, Lars von den Driesch, Clovis Galiez, Maria J. Martin, Johannes Söding, and Martin 649 Steinegger. Uniclust databases of clustered and deeply annotated protein sequences and alignments. 650 Nucleic Acids Research, 45(D1):D170–D176, 2017. ISSN 0305-1048. doi: 10.1093/nar/gkw1081. 651 Milot Mirdita, Konstantin Schütze, Yoshitaka Moriwaki, Lim Heo, Sergey Ovchinnikov, and Martin 652 Steinegger. ColabFold: making protein folding accessible to all. Nat Methods, 19(6):679–682, 653 2022. doi: 10.1038/s41592-022-01488-1. 654 655 Eric Nguyen, Michael Poli, Marjan Faizi, Armin W Thomas, Michael Wornow, Callum Birch-Sykes, 656 Stefano Massaroli, Aman Patel, Clayton M. Rabideau, Yoshua Bengio, Stefano Ermon, Christopher 657 Re, and Stephen Baccus. HyenaDNA: Long-range genomic sequence modeling at single nucleotide 658 resolution. In Thirty-seventh Conference on Neural Information Processing Systems, 2023. URL 659 https://openreview.net/forum?id=ubzNoJjOKj. 660 661 Eric Nguyen, Michael Poli, Matthew G. Durrant, Armin W. Thomas, Brian Kang, Jeremy Sullivan, Madelena Y. Ng, Ashley Lewis, Aman Patel, Aaron Lou, Stefano Ermon, Stephen A. Baccus, Tina 662 Hernandez-Boussard, Christopher Ré, Patrick D. Hsu, and Brian L. Hie. Sequence modeling and 663 design from molecular to genome scale with Evo. bioRxiv, 2024. doi: 10.1101/2024.02.27.582234. 664 URL https://www.biorxiv.org/content/10.1101/2024.02.27.582234v1. 665 666 Erik Nijkamp, Jeffrey A Ruffolo, Eli N Weinstein, Nikhil Naik, and Ali Madani. Progen2: exploring 667 the boundaries of protein language models. *Cell systems*, 14(11):968–978, 2023. 668 669 Pascal Notin, Lood Van Niekerk, Aaron W Kollasch, Daniel Ritter, Yarin Gal, and Debora S. Marks. 670 TranceptEVE: Combining family-specific and family-agnostic models of protein sequences for 671 improved fitness prediction. *bioRxiv*, 2022. doi: 10.1101/2022.12.07.519495. URL https: //www.biorxiv.org/content/early/2022/12/27/2022.12.07.519495. 672 673 Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood van Niekerk, Steffanie Paul, Han Spinner, Nathan 674 Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, Jonathan Frazer, Mafalda Dias, Dinko 675 Franceschi, Yarin Gal, and Debora Marks. ProteinGym: Large-scale benchmarks for protein fitness 676 prediction and design. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine 677 (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 64331–64379. Curran 678 Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/ 679 paper/2023/file/cac723e5ff29f65e3fcbb0739ae91bee-Paper-Datasets_ 680 and_Benchmarks.pdf. 681 Zhangzhi Peng, Benjamin Schussheim, and Pranam Chatterjee. PTM-Mamba: A PTM-aware 682 protein language model with bidirectional gated Mamba blocks. bioRxiv, 2024. doi: 10.1101/ 683 2024.02.28.581983. URL https://www.biorxiv.org/content/early/2024/02/ 684 29/2024.02.28.581983. 685 686 Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y Fu, Tri Dao, Stephen Baccus, Yoshua 687 Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional 688 language models. In International Conference on Machine Learning, pp. 28043–28078. PMLR, 689 2023. 690 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi 691 Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text 692 transformer, 2020. URL https://arxiv.org/abs/1910.10683. 693 694 Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu, and Alexander Rives. MSA Transformer. In Proceedings of the 38th International Conference on 696 Machine Learning, volume 139, pp. 8844–8856. PMLR, 2021. URL https://proceedings. 697 mlr.press/v139/rao21a.html. Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, 699 Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function 700 emerge from scaling unsupervised learning to 250 million protein sequences. Proc. Natl. Acad. 701 Sci. U.S.A., 118(15), 2021. ISSN 0027-8424. doi: 10.1073/pnas.2016239118.

702 703 704 705	William P. Russ, Matteo Figliuzzi, Christian Stocker, Pierre Barrat-Charlaix, Michael Socolich, Peter Kast, Donald Hilvert, Remi Monasson, Simona Cocco, Martin Weigt, and Rama Ranganathan. An evolution-based model for designing chorismate mutase enzymes. <i>Science</i> , 369(6502):440–445, 2020. ISSN 10959203. doi: 10.1126/science.aba3304.
706 707 708	D. Sgarbossa, U. Lupo, and AF. Bitbol. Generative power of a protein language model trained on multiple sequence alignments. <i>Elife</i> , 12:e79854, 2023. doi: 10.7554/eLife.79854.
709 710 711	Jin Su, Chenchen Han, Yuyang Zhou, Junjie Shan, Xibin Zhou, and Fajie Yuan. SaProt: Protein language modeling with structure-aware vocabulary. <i>bioRxiv</i> , 2023. doi: 10.1101/2023.10.01. 560349.
712 713 714	The UniProt Consortium. UniProt: the universal protein knowledgebase in 2021. <i>Nucleic Acids Research</i> , 49(D1):D480–D489, 2021. ISSN 0305-1048. doi: 10.1093/nar/gkaa1100.
715 716 717	J. D. Thompson, B. Linard, O. Lecompte, and O. Poch. A comprehensive benchmark study of multiple sequence alignment methods: current challenges and future perspectives. <i>PLoS One</i> , 6(3): e18093, Mar 2011.
718 719 720 721 722	Timothy Truong Jr and Tristan Bepler. PoET: A high-performing protein lan- guage model for zero-shot prediction. https://www.openprotein.ai/ poet-a-high-performing-protein-language-model-for-zero-shot-prediction. Accessed: 2024-05-21.
723 724	Timothy Truong Jr and Tristan Bepler. PoET: A generative model of protein families as sequences- of-sequences. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
725 726	M. van Kempen, S. S. Kim, C. Tumescheit, M. Mirdita, J. Lee, C. L. M. Gilchrist, J. Söding, and M. Steinegger. Fast and accurate protein structure search with Foldseek. <i>Nat. Biotechnol.</i> , 2023.
728 729 730 731	Robert Verkuil, Ori Kabeli, Yilun Du, Basile I. M. Wicky, Lukas F. Milles, Justas Dauparas, David Baker, Sergey Ovchinnikov, Tom Sercu, and Alexander Rives. Language models generalize beyond natural proteins. <i>bioRxiv</i> , 2022. doi: 10.1101/2022.12.21.521521. URL https://www.biorxiv.org/content/early/2022/12/22/2022.12.21.521521.
732 733 734 735	Jesse Vig, Ali Madani, Lav R. Varshney, Caiming Xiong, Richard Socher, and Nazneen Rajani. BERTology meets biology: Interpreting attention in protein language models. In <i>International Conference on Learning Representations</i> , 2021. URL https://openreview.net/forum? id=YWtLZvLmud7.
736 737 738 739 740	Roger Waleffe, Wonmin Byeon, Duncan Riach, Brandon Norick, Vijay Korthikanti, Tri Dao, Albert Gu, Ali Hatamizadeh, Sudhakar Singh, Deepak Narayanan, Garvit Kulshreshtha, Vartika Singh, Jared Casper, Jan Kautz, Mohammad Shoeybi, and Bryan Catanzaro. An empirical study of mamba-based language models. <i>arXiv</i> , pp. 2406.07887, 2024. URL https://arxiv.org/abs/2406.07887.
742 743	K. K. Yang, N. Fusi, and A. X. Lu. Convolutions are competitive with transformers for protein sequence pretraining. <i>Cell Syst</i> , 15(3):286–294, Mar 2024.
744	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	

760

809

Supplementary material

A PROTEINGYM ASSAYS USED IN VALIDATION

Here, we list the 20 assays we extracted from the ProteinGym benchmark to choose some hyperparameters (see Figure S8).

764	$4 \qquad \qquad$	Melnikov 2014					
765	AMER HUMAN Tsuboyama 2023 4G30 PITX2 HUMAN	Tsuboyama 2023 2L7M					
766	6 CAR11 HUMAN Meitlis 2020 lof PPM1D HUMAN	N Miller 2022					
767	7 CBS HUMAN Sun 2020 R1AB SARS2 F	lynn 2022					
768	⁸ CUE1 YEAST Tsuboyama 2023 2MYX RDRP I33A0 Li	2023					
769	⁹ DYR_ECOLI_Nguyen_2023 S22A1_HUMAN	Yee_2023_abundance					
770	⁰ GDIA_HUMAN_Silverstein_2021 SCN5A_HUMAN	N_Glazer_2019					
771	¹ HIS7_YEAST_Pokusaeva_2019 SHOC2_HUMAN	N_Kwon_2022					
772	2 HXK4_HUMAN_Gersing_2023_abundance TRPC_SACS2_C	han_2017					
773	3 KCNE1_HUMAN_Muhammad_2023_expr VILI_CHICK_Ts	uboyama_2023_1YU5					
774	4						
775	5						
776 777	⁶ B DETAILS ON MUTATIONAL EFFECT PREDICTION W	ITH PROTMAMBA					
778	⁸ In Table 1, we consider different ways of computing mutational e	effects using ProtMamba and					
779	⁹ compare with other models. Here, we explain our different approach	es in more detail.					
780							
781	• The entry ProtMamba (single) of Table 1 is the result	reported when using the FIM					
782	2 technique at the end of the sequence to evaluate mutations.	It is the fastest method among					
783	3 homology-aware ones. The procedure is the following:						
784	4 1. Subsample a predetermined number of homologs of the	e target sequence considered in					
785	5 the DMS to be used as context, based on a diversity filt	ter.					
786	6 2. Run ProtMamba on the context and collect the last hid	den state of the context (as the					
787	7 model is recurrent).	model is recurrent).					
788	3. Start from the last hidden state as initial state for every v	variant to score. Scoring is then					
789	⁹ very fast, since we only need to apply Prolivianida on and not on the full context. As Membe scales linearly is	a single sequence per variant,					
790	⁰ evaluate many different variants very fast on a single G	PU (hundreds to thousands per					
791	¹ batch). The mutated residues are put at the end of the s	sequence in the FIM mask. For					
792	single mutations, we can then evaluate in one shot the likelihood of every mutation.						
793	4. We compare the likelihood of the WT to the likelihood of the variant using Fill-in-the-						
794	middle to get a proxy of variant fitness. We can then evaluate the fitness of variant						
795	5 comprising each of the 20 amino-acids at the mutated s	site(s) in a single shot.					
796	• The entry ProtMamba AR (autoregressive) of Table 1 is the	result reported when evaluating					
797	7 the likelihood of a variant without using the FIM technique.	The procedure is:					
798	⁸ 1. Perform steps 1 and 2 as in ProtMamba (single) above.	1					
799	9 2. Start from the last hidden state as initial state for every va	riant to score. We then evaluate					
800	0 the autoregressive likelihood of the full variant. Since	some logits are computed after					
801	1 the mutation in this setup (because they are positioned a	fter in the sequence), we cannot					
802	2 evaluate the fitness of the 20 amino-acids in a single s	hot (in contrast to the previous					
803	approach), which increases the number of calls to Pr	otMamba, and hence requires					
804	4 more time.						
805	5 3. Compare the likelihood of the WT over the full sequence	e to the likelihood of the variant					
806	6 to get an evaluation of variant fitness.						
807	• The entry ProtMamba (w/ R), i.e. with retrieval of Table	e 1 is the result reported when					
808	s combining ProtMamba with a prior based on the frequency	of amino-acids in the MSA of					

- the relevant protein family.
 - 1. Perform steps 1, 2 and 3 as in ProtMamba (single) above.

2. Load the MSA, compute a log-likelihood prior (log $p_{retrieval}$) based on the frequency of every amino-acid, and sum it with ProtMamba's log-likelihood (log $p_{ProtMamba}$):

 $\log p_{\text{ProtMamba}(w/R)} = \alpha \log p_{\text{retrieval}} + (1 - \alpha) \log p_{\text{ProtMamba}}.$

814The parameter α was optimized using the validation set introduced in Section A (see815Figure S8). This operation can be parallelized across CPUs or GPUs. We report the816results using 16 workers.817

Note that we tested ProtMamba AR both with our model fine-tuned on FIM (ProtMamba Long
Finetuned) and with our foundation model (ProtMamba Long), and we obtained similar results
(respectively 0.361 and 0.367).

In Table S1, we break down results by MSA depth and by number of mutations. We observe that for
datasets with more than one mutation (last column in Table S1), ProtMamba with retrieval slightly
outperforms the overall state-of-the-art model TranceptEVE L and reaches performance close to the
structure-based model ESM-IF1. However, averaging over all datasets, ProtMamba does not reach
the same performance as TranceptEVE L. But since ProtMamba performs better than Tranception L,
ensembling ProtMamba and EVE predictions might yield comparable performance.

		Spearman correlation by MSA depth				by mutations	
Model	Par.	All depths	Deep	Medium	Shallow	1	2+
ESM-2	150M	0.387	0.497	0.358	0.306	0.367	0.379
ESM-IF1	142M	0.422	0.544	0.431	0.300	0.413	0.471
Tranception S (w/o R)	85M	0.303	0.320	0.295	0.258	0.293	0.262
Tranception L (w/o R)	700M	0.374	0.419	0.371	0.358	0.358	0.390
ProtMamba (w/o R)	107M	0.406	0.465	0.411	0.391	0.376	0.444
MSA Transformer	100M	0.421	0.473	0.435	0.393	0.392	0.435
Tranception S (w/ R)	85M	0.418	0.444	0.415	0.428	0.389	0.409
Tranception L (w/ R)	700M	0.434	0.473	0.438	0.432	0.404	0.463
TranceptEVE L	>700M	0.456	0.492	0.467	0.451	0.426	0.467
ProtMamba (w/ R)	107M	0.432	0.472	0.438	0.448	0.404	0.469

Table S1: **Performance of different models on the ProteinGym benchmark.** We report Spearman correlation values obtained both based on retrieval (w/ R) and non-retrieval (w/o R) methods, and parameter count for each model. We report results divided according to MSA depth and number of mutations in the benchmark dataset. Results for benchmark models were obtained from https://proteingym.org/. Note that PoET-205M (Truong Jr & Bepler, 2024) reports an overall Spearman correlation of 0.474 (Truong Jr & Bepler) on ProteinGym, but it is not yet on the ProteinGym website, and no information is given about the training time or resources.

C ABLATION STUDIES

We investigated ablations or alternative implementations of ProtMamba using two models: a small
model with 14 million parameters (8 layers, hidden dimension 512) and the standard architecture
with 107 million parameters (16 layers, hidden dimension 1024). Both models were trained for 10
billion tokens (50k steps with a batch size of 128) and evaluated on a validation set of 500 unseen
clusters. During evaluation, a context of 25 sequences was used. The perplexity of the models was
assessed both autoregressively (left-to-right) and in the fill-in-the-middle (FIM) spans.

The performance of these alternatives (described in more detail below) is summarized in Table S2.
Perplexity values are reported for both autoregressive and FIM modes in the small model and the larger one.

864 865	Pernlevity	14M Para	ameters	107M Parameters		
866	respicately	Autoregressive FIM		Autoregressive	FIM	
867	Only FIM from scratch	Fail	13.90 ± 0.34	Fail	15.59 ± 0.27	
868	AR only	12.58 ± 0.31	18.03 ± 0.25	11.05 ± 0.36	Fail	
869	No positional encoding	13.01 ± 0.30	16.71 ± 0.47	12.31 ± 0.37	17.20 ± 0.58	
870	Additive positional encoding	12.72 ± 0.31	13.60 ± 0.33	12.58 ± 0.38	13.81 ± 0.31	
871	One mask, one token	12.76 ± 0.31	15.54 ± 0.29	11.04 ± 0.33	16.60 ± 0.36	
872	Masking fraction 50%	13.02 ± 0.31	$\textbf{13.44} \pm 0.33$	$\textbf{10.94} \pm 0.36$	$\textbf{11.59} \pm 0.35$	
873	ProtMamba	13.00 ± 0.30	13.89 ± 0.32	11.35 ± 0.33	12.62 ± 0.30	
874						

Table S2: Alternatives to ProtMamba. Perplexity values are reported for different alternatives to ProtMamba, evaluated on small (14M parameters) and larger (107M parameters) models, for both autoregressive and FIM tasks.

The alternative implementations tested against our main ProtMamba model, and whose performance is reported in Table S2, were constructed as follows:

- **Only FIM from scratch:** This approach backpropagates the loss exclusively from the FIM tokens, disregarding the main amino-acid chain. Training this way from scratch disables autoregressive (left-to-right) next-token prediction and degrades performance, including on FIM tasks.
 - Autoregressive (AR) only: Trains the model without sampling FIM spans. While this slightly improves autoregressive performance, it significantly degrades FIM capabilities.
 - **No positional encoding:** Omits positional encodings entirely. In autoregressive mode, the model can partially rely on its recurrent architecture, but in FIM mode, performance suffers due to the absence of positional information in input.
- Additive positional encoding: Uses additive positional encoding (summing token embeddings with positional encoding) instead of concatenated positional encoding (concatenating token embeddings with positional encoding). This approach showed mixed results, with slight improvement in the small model but degradation in the larger model.
- One mask, one token: Uses one mask per token (as in the T5 model (Raffel et al., 2020)) instead of one mask per span of tokens (as in our approach, inspired by Bavarian et al. (2022)). This approach led to performance degradation in FIM, likely due to insufficient training on larger number of mask tokens.
 - **Masking fraction 50%:** Samples 50% of the tokens for FIM (compared to ProtMamba's 20%). This alternative brought minor but noticeable improvements, suggesting potential for further development.

918 D SUPPLEMENTARY FIGURES



Figure S1: Loss and perplexity during training. Cross entropy loss and perplexity computed for both the full non-masked sequences and the FIM tokens. We show them as a function of the number of tokens processed during the training of ProtMamba. They are computed on the training set and on a validation set of 192 held-out OpenProteinSet sequence clusters (see Section 2.3).







Figure S4: Loss and perplexity of the full sequences vs. number of sequences in the context. Scaling of the per-sequence perplexity (i.e. the standard autoregressive perplexity of the full nonmasked sequence) versus the number of context sequences. Results are averaged over all 500 clusters of the test set and 20 replicates for each cluster (differing by the random sampling of context sequences). Context sizes go up to 2^{17} amino acids. Sequence clusters are split according to the average length *L* of sequences in the cluster. We observe that clusters with shorter sequences reach lower perplexities.

- 1023 1024
- 1025





Figure S7: Impact of context length on results on the ProteinGym benchmark. (a) We run
ProtMamba Long on the ProteinGym dataset, building contexts of different sizes in terms of numbers
of tokens (from 8,000 to 128,000). We see that the increase in performance is more important for
long sequences, which highlights the benefit of long context to model long protein sequences. (b) We
also run ProtMamba Long on the ProteinGym dataset, building contexts of different sizes in terms
of numbers of sequences (from 25 to 200). Overall, we notice a rise in the Spearman correlation,
showing that prediction benefits from longer context.



1124 Figure S8: Choice of context length and retrieval coefficient using a validation set. We randomly 1125 extracted a validation set of 20 datasets (see supplementary Section A) to select the best context length and retrieval coefficient. (a) The prediction improves with the context size in the validation set. 1126 This trend was later observed in the rest of the benchmark (testing set) too. (b) Retrieval requires 1127 mixing the fitness score \mathcal{F}_m obtained from ProtMamba and the fitness score obtained from the 1128 independent-site model \mathcal{F}_i through the retrieval fitness score $\mathcal{F}_r = \alpha \mathcal{F}_i + (1 - \alpha) \mathcal{F}_m$. The best 1129 model on the validation set was obtained for a retrieval coefficient $\alpha = 0.5$, which was later verified 1130 on the rest of the dataset. 1131

1104 1105 1106



Figure S10: Example results on the ProteinGym benchmark. Results of ProtMamba Long are
shown on 25 randomly sampled deep mutational scan (DMS) experimental datasets from ProteinGym,
and are compared to existing methods (see main text). The score shown is the Spearman correlation
between predictions and experimental results.



variants in context (b) or using active and inactive variants in context (c). Inactive variants tend to have higher perplexity (implying lower fitness score) when the context contains only active variants (b) while active variants have lower perplexity (implying higher fitness score) when the context contains only active variants (c).





Figure S14: **ProtMamba perplexity versus DCA energy for chorismate mutase variants.** Prot-Mamba perplexity is evaluated using full sequences, FIM and only active variants in the context, and is shown versus the Potts or DCA energy from Russ et al. (2020). Active variants are in green, while inactive variants are in red. We observe that most of the variants that are active have low perplexity, and that many inactive variants that were not discriminated as inactive by DCA are labelled as such by ProtMamba (bottom right part of the plot).

1328	Cluster	Hamming	HMMER	pLDDT	pTM	F2CV06 -	Hamming
1330	A0A2H9MP70	0.45	-0.44	-0.77	-0.54	A0A1C5UJ41 - A0A1A8YWK1 -	🔺 HMMER 🔒 💏
1331	G4ZH78	0.23	-0.08 -0.42	-0.01 -0.47	-0.38 -0.47	A0A1S3G530 - D8SD16 -	▼ pLDDT ▲★ ■
1332 1333	A0A0A0HZM8	0.79	-0.7 -0.47	-0.84 -0.69	-0.81 -0.63	A0A146ZGL6 -	★ Mean
1334	A0A0911DH7 A0A2N1P554	0.31	-0.32 -0.35 0.79	-0.46 -0.37	-0.46 -0.4	A0A194V424 - S7UZ45 -	• ***
1335 1336	A0A1C5UJ41 A0A194V424	0.81	-0.78 -0.71 0.57	-0.81 -0.7 0.74	-0.77 -0.54	A0A0A0HZM8 -	
1337	F2CV06	0.50	-0.57 -0.85 0.75	-0.74 -0.81 0.56	-0.00 -0.79	A0A241VGM5 -	
1339	D8SD16	0.65	-0.73 -0.74	-0.50 -0.71 0.15	-0.61	A0A2X4BAY2 -	
1340	A0A159IN77 A0A1C6Q5J2	0.18	-0.22 -0.22	-0.15 -0.45	-0.27 -0.27	A0A091TDH7 -	
1342	A0A2X4BAY2	0.44 0.27	-0.52 -0.63	-0.61 -0.63	-0.52 -0.54	A0A12000352	****
1343 1344	A0A1S3G530 Mean	$0.44 \\ 0.88 \\ 0.54$	$-0.59 \\ -0.72 \\ -0.56$	-0.62 -0.8 -0.62	$-0.55 \\ -0.74 \\ -0.57$		0.2 0.4 0.6 0.8
1045						-	realson correlation

Figure S15: **Pearson correlation between ProtMamba perplexity and scores for generated sequences.** For each of 19 test clusters, we used all the sequences generated by ProtMamba to compute the Pearson correlation between the model perplexity and the Hamming distance to the closest natural neighbor, the HMMER score, the pLDDT and pTM scores from ESMFold.





Figure S16: **Properties of generated sequences.** Left panels: histograms of Hamming distances, HMMER scores and mean pLDDT scores from ESMFold of generated sequences for 10 example test clusters (10 rows). Right panels: scatter plots of ProtMamba perplexity versus the Hamming distance to the closest natural neighbor, the HMMER score and mean pLDDT score from ESMFold for all generated sequences from each of 10 example clusters (10 rows). Dashed vertical lines: median of the generated sequences (blue), median of the natural sequences (green) and pLDDT value of the reference structure of the cluster (red). The last one is shown only for the rightmost plot.