

---

# FLOP: Tasks for Fitness Landscapes Of Protein wildtypes

---

**Peter Mørch Groth<sup>1,2</sup>, Richard Michael<sup>1</sup>,  
Jesper Salomon<sup>2</sup>, Pengfei Tian<sup>2</sup>, Wouter Boomsma<sup>1</sup>**

<sup>1</sup>Department of Computer Science, University of Copenhagen

<sup>2</sup>Bioinformatics & Design, Enzyme Research, Novozymes

{petergroth, richard.michael, wb}@di.ku.dk

{pmg, jrsox, pfi}@novozymes.com

## Abstract

1 Protein engineering has the potential to create optimized protein variants with  
2 improved properties and function. An initial step in the protein optimization pro-  
3 cess typically consists of a search among natural (wildtype) sequences to find  
4 the naturally occurring proteins with the most desirable properties. Promising  
5 candidates from this initial discovery phase then form the basis of the second step:  
6 a more local optimization procedure, exploring the space of variants separated  
7 from this candidate by a number of mutations. While considerable progress has  
8 been made on evaluating machine learning methods on single protein datasets,  
9 benchmarks of data-driven approaches for global fitness landscape exploration are  
10 still lacking. In this paper, we have carefully curated a representative benchmark  
11 dataset, which reflects industrially relevant scenarios for the initial wildtype discov-  
12 ery phase of protein engineering. We focus on exploration within a protein family,  
13 and investigate the downstream predictive power of various protein representation  
14 paradigms, i.e., protein language model-based representations, structure-based  
15 representations, and evolution-based representations. Our benchmark highlights  
16 the importance of coherent split strategies, and how we can be misled into overly  
17 optimistic estimates of the state of the field. The codebase and data can be accessed  
18 via <https://github.com/petergroth/FLOP>.

## 19 1 Introduction

20 The goal of protein engineering is to optimize proteins towards a particular trait of interest. This has  
21 applications both for industrial purposes and drug design. There is clear potential for machine learning  
22 to aid in this process. By predicting which protein sequences are most promising for experimental  
23 characterization, we can accelerate the exploration of the “fitness landscape” of the protein in question  
24 [1]. Regression of functional landscapes is challenging for multiple reasons. Typically a data-scarce  
25 problem, careful considerations of the experimental setup are required to avoid inadvertent data  
26 leakage. Concerning the functional landscapes of naturally occurring (also known as *wildtype*)  
27 proteins, the pairwise amino acid sequence identities can often vary significantly with some proteins  
28 differing by only a single amino acid while others might be less than ten percent similar. High-  
29 throughput experimental techniques are improving the data scarcity issue, while underlying structure  
30 typically exists in the datasets allowing for supervised learning despite the intrinsic challenges.

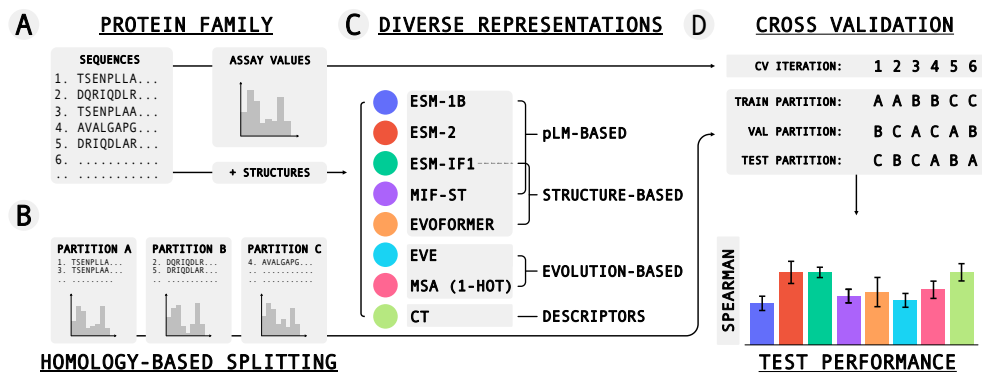


Figure 1: Schematic over dataset splitting, representations, and cross-validation process. A: A dataset with sequences from a single protein family and corresponding assay values is curated. B: A stratified sequence identity splitting procedure generates partitions A, B, and C, which are (1) homologically different from each other, (2) contain similar number of sequences, and (3) match the full dataset’s target distribution. C: Eight types of protein representations are computed. D: Cross-validation using a random forest regressor is applied to obtain mean values and standard errors on the test partitions.

31 The optimization process of proteins and enzymes can typically be divided into multiple stages. An  
 32 often employed initial step is to search for promising candidates among wildtype proteins, resulting  
 33 in a set of proteins with desirable properties. We will refer to this as the *wildtype discovery* phase [2].  
 34 Since we are typically optimizing for a specific trait, we may often limit this initial exploration to a  
 35 particular protein family, where the members share an evolutionary history which has resulted in a  
 36 similar function. The selected set of wildtype proteins will then form the basis for a second phase  
 37 in the engineering process: localized optimization, where novel variants of the wildtype proteins  
 38 are examined through various assays [3]. Sometimes, the wildtype discovery phase is not only  
 39 carried out once as several rounds might be required before an initial suitable candidate is found.  
 40 Additionally, the resulting candidate might prove insufficient at a later stage of protein engineering,  
 41 where conditions such as temperature are altered or where stress-factors are introduced.

42 In recent years, we have seen considerable efforts in defining benchmarks to help the machine learning  
 43 community make progress in this field. However, these efforts have primarily focused on the second  
 44 stage, i.e., variant effect prediction, where a dataset consist of thousands of variants from a single  
 45 wildtype. In this paper, we argue for the importance of establishing well-defined benchmark tasks for  
 46 the first stage as well. We present three challenging tasks and a careful analysis of the experimental  
 47 design, demonstrating how poor choices can lead to dramatic overestimation of performance.

48 We conduct our experiments using a variety of fixed-size protein representations: sequence-based  
 49 embeddings obtained through protein language models, structure-based representations from fold-  
 50 ing and inverse folding models, evolution-based representations obtained from multiple sequence  
 51 alignments, as well as simple biologically-motivated sequence descriptors. In addition to the su-  
 52 pervised approach, we include four zero-shot predictors to showcase a simpler approach to the task  
 53 of identifying promising candidates. We show that the choice of representation can greatly affect  
 54 the downstream predictive performance, and we therefore argue that more progress can be made by  
 55 constructing meaningful representations and not solely in the construction of complex prediction  
 56 models. Given the oftentimes limited dataset sizes, we therefore rely on a random forest regressor.

## 57 2 Related work

58 Benchmarks play an important role in driving progress in protein-related prediction tasks. The  
 59 most well-known is perhaps the rolling CASP benchmark, which is arguably responsible for the  
 60 recent breakthroughs in protein structure prediction [4–6]. For the prediction of protein stability and  
 61 function, several studies have curated relevant experimental datasets for use as benchmarks. The  
 62 TAPE benchmark was an early such example designed to test protein sequence representations on a

63 set of diverse downstream tasks [7]. Two of these tasks were related to protein engineering: stability  
64 prediction on variants of a set of 12 designed proteins [8] and characterization of the functional  
65 landscape of green fluorescent protein [9]. The PEER benchmark [10] expanded on the TAPE  
66 benchmark with many additional tasks. This included prediction of  $\beta$ -lactamase activity [11], and a  
67 binary solubility classification task on a diverse set of proteins. Focusing entirely on variant effects,  
68 the recent ProteinGym benchmark has assembled a large set of Deep Mutational Scanning (DMS)  
69 assays and made them available as substitution and insertion-deletion prediction tasks [12]. While  
70 the above all consider protein sequence inputs, the recent Atom3D benchmark [13] presents various  
71 prediction tasks using 3D structure as input, including predicting amino acid identity from structural  
72 environments (for general proteins), and mutation effects on protein binding, using data originating  
73 from the SKEMPI database [14, 15].

74 Most closely related to this current paper is the FLIP benchmark, which dedicates itself to the  
75 prediction of functional fitness landscapes of proteins for protein engineering [16]. FLIP introduces  
76 three tasks: one on the prediction of protein stability of wildtype proteins (distributed over many  
77 families) using data from the Meltome Atlas [17], and two tasks focused on mutations at specific  
78 sites of proteins GB1 [18] and AAV [19]. While the FLIP benchmark is of great value for protein  
79 engineering, there are key characteristics which make it unsuitable for wildtype discovery, e.g.,  
80 the use of the Meltome Atlas, which consists of thousands of sequences from different organisms  
81 spanning many different protein families. The sequences in the GB1 dataset only have mutations at  
82 four fixed positions while the sequences in the AAV dataset only contain 39 mutation sites, both of  
83 which corresponds to mutations at less than 10% of the full-length proteins. Such datasets with very  
84 local fitness landscapes are not generalizable enough for wildtype discovery.

85 Most functional tasks in current benchmarks are thus concerned with protein sequences that are  
86 derived from a single wildtype sequence by one or more mutations. Characterizing the functional  
87 effects of such variants is critical for protein engineering. However, before engaging in the optimiza-  
88 tion process itself, it is important to select meaningful starting points. As a natural complement to the  
89 FLIP benchmark, we therefore present a novel benchmark titled FLOP. The tasks we present are the  
90 characterization of functional landscapes of wildtype proteins.

91 Our curated datasets all consist of functionally characterized wildtype sequences. For each dataset, we  
92 limit ourselves to a single family, and define our tasks as regression problems on the functional assay  
93 values. While mutational fitness landscape datasets are relatively abundant, few published datasets  
94 exist where the global fitness landscapes of wildtype proteins from single families are examined. This  
95 imposes limitations in the number and sizes of available datasets which are suitable for our considered  
96 problem. Given the low-data regime, the focus of our benchmark is thus to find representations of the  
97 protein input that makes few-shot or even zero-shot learning feasible. As a point of departure, we  
98 provide a set of state-of-the-art embeddings, reflecting different protein modalities.

### 99 **3 Experimental setup**

100 The domain we explore in this work is characterized by data scarcity, requiring special care in the  
101 design of the experimental setup. Figure 1 shows an overall schematic of the benchmarking process.

#### 102 **3.1 Dataset splitting**

103 With the proliferation of large datasets and computationally demanding models, a common learning  
104 paradigm in machine learning is to rely on hold-out validation, whereby fixed training, validation,  
105 and testing sets are randomly generated. This method has several serious limitations when applied to  
106 biological datasets of limited sizes. Firstly, randomly splitting a dataset assumes that the data points  
107 are independent and identically distributed (i.i.d.). This is however not the case for members of a  
108 protein family which share common ancestors, leading to potential data leakage if protein sequences  
109 that are close in evolutionary space are placed in separate splits. Secondly, when splitting small  
110 datasets for a hold-out validation approach for supervised learning, the target values might not be  
111 well-balanced, resulting in dissimilar target distributions thus leading to bias and poor generalizability.

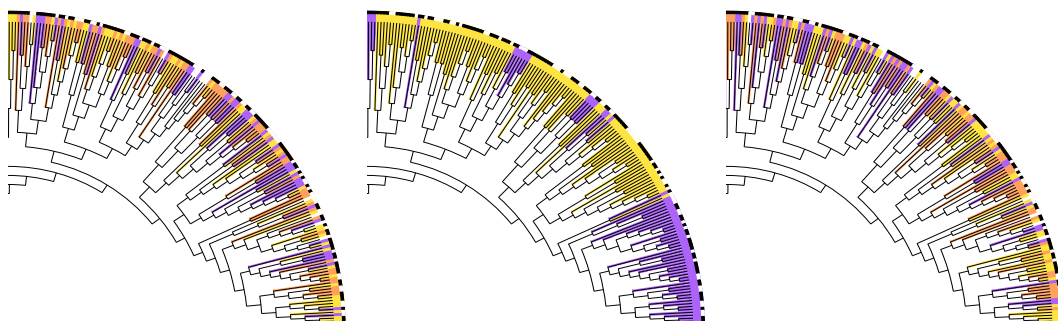


Figure 2: The same segment of a phylogenetic tree for the PPAT dataset. Branch color corresponds to its CV partition, while the outermost ring shows the target values (black indicates high and white indicates low values). The segments highlight the diversity found in wildtype protein families. Left: entries are colored according to the prescribed dataset splitting procedure which allows learning across subfamilies (indicated by the mix of colors). Middle: entries are colored by a clustering approach leading to wide regions, inhibiting learning across subfamilies. Right: entries are randomly assigned a color. While similar to the leftmost scheme, the random coloring allows near identical sequences to be placed in separate partitions leading to excessive data-leakage.

112 To handle these potential issues, we rely on a sequence identity-based, stratified cross-validation  
 113 procedure ensuring that (1) partitions are generated such that any two proteins occurring in different  
 114 partitions are guaranteed to be different at a pre-set homology cut-off, (2) cross-validation (CV)  
 115 minimizes the potential bias which might occur during hold-out validation, (3) the target distribution  
 116 is reflected by the generated partitions via stratification on discretized target values, and (4) the  
 117 number of sequences in each partition is similar to reduce the variance.

118 To generate these high-quality data partitions, we use the four-phase procedure described in [20]  
 119 and implemented in the GraphPart framework [21] to create three label-balanced partitions for  
 120 each dataset. We begin the procedure from an initial sequence identity threshold, and increase the  
 121 threshold until the generated partitions are of sufficient sizes (i.e., at least 25 % of sequences in all  
 122 three partitions). The stratification is achieved by creating a binary label which indicates whether a  
 123 protein has low or high target value, e.g., by fitting a two-component Gaussian mixture model. For  
 124 the dataset-specific stratification boundaries, see Section A in the supplementary materials.

125 Figure 2 shows the same segment of a phylogenetic tree of the curated PPAT dataset, showing the  
 126 evolutionary relationship between sequences. Large versions of the trees can be found in Section E.  
 127 The colors indicate which CV partition each sequence belongs to while the black and white squares  
 128 in the outer ring indicate the stratification labels. The segments show the diversity encountered in  
 129 wildtype protein families. The left segment is colored by our splitting procedure and shows that it  
 130 manages to create diverse partitions spanning the entire evolutionary tree to allow learning across  
 131 protein subfamilies. The middle segment is colored by an MMseqs [22] clustering approach, leading  
 132 to contiguous areas inhibiting learning across subfamilies. The entries in the rightmost segment are  
 133 randomly assigned a color, corresponding to random splitting. While similar to the leftmost scheme,  
 134 the random coloring allows near-identical sequences to be placed in separate partitions leading to  
 135 excessive data-leakage.

### 136 3.2 Representations

137 To accurately reflect the current paradigms of state-of-the-art protein representations, we choose  
 138 representatives from three main categories, the dimensionalities of which can be found in Section G  
 139 in the supplementary materials.

140 Protein language models (pLMs) that are trained on hundreds of millions of protein sequences in an  
 141 unsupervised fashion have been proven to be competitive for a multitude of tasks including supervised  
 142 prediction of protein properties, residue contact prediction, variant effect prediction [16, 23–26],  
 143 etc. We here choose the popular ESM-1B [24] and the more recent ESM-2 models [27]. To fix the

144 dimensionality for proteins of different lengths, we perform mean-pooling over the residue dimension.  
145 This operation is likely to filter out information encoded along the protein sequence and more optimal  
146 approaches will likely yield more informative representations and thus higher predictive performance  
147 (see Table 2 in [28]). Constructing fixed-size embeddings from sequences of variable lengths is  
148 however nontrivial and considered out of the scope of this study.

149 The second category we include is structure-based. We extract embeddings from the Evoformer-  
150 modules while folding proteins with AlphaFold2 [29] via ColabFold [30], which have been shown  
151 to perform well for structure-related prediction tasks [31]. Using the predicted structures, we then  
152 extract embeddings from the inverse-folding model ESM-IF1 (also known as the GVP-GNN) [32]  
153 which incorporates a pLM and graph neural network architecture. We similarly use embeddings from  
154 the MIF-ST model, which is an inverse folding model leveraging a pretrained convolutional pLM  
155 [33]. As with the pLMs, we apply mean pooling to achieve sequence-level embeddings.

156 The third category is evolution-based. As a baseline, we will use a one-hot encoded multiple  
157 sequence alignment (MSA) over the proteins of interest [34–36]. Since the MSA is independent  
158 of labels, we enrich the unaligned sequence pools with additional members from the respective  
159 protein families using UniProt [37] and InterPro [38]. Given MSAs, models can be designed which  
160 leverage the evolutionary history of the protein family (e.g., EVE and related models [12, 39–42]).  
161 For each curated dataset, we train EVE [39] on the corresponding protein family and extract the latent  
162 representations. Technical details on the training procedure can be found in Section H.

163 In addition to these groups of advanced representations, we include compositional and transitional  
164 (CT) physicochemical descriptors for each protein sequence as a simple baseline, which relate to  
165 overall polarizability, charge, hydrophobicity, polarity, secondary structure, solvent accessibility, and  
166 van der Waals volume of each sequence as predicted using the PyBioMed library [43].

167 With the exception of the physicochemical descriptors, all included representations rely on models  
168 which have been pretrained on thousands to hundreds of millions of proteins. While it is possible that  
169 a number of the sequences in the curated datasets also belong to the training sets of these models  
170 (which by design is the case for the evolution-based approaches), we do not consider this to be a  
171 fatal form of data leakage as it purely pertains to the un- or self-supervised pretraining phases and is  
172 independent of the sequence labels.

### 173 **3.3 Regression**

174 The purpose of this benchmark is to provide a structured procedure to evaluate the predictive  
175 performance on downstream regression tasks given protein representations. We believe that larger  
176 prediction improvements can be achieved by focusing on developing novel protein representations  
177 rather than more complex regression models. Due to the low-N setting in which we operate, the  
178 training of large, complex models is practically inhibited, which is why we have chosen to rely  
179 on a random forest regressor. For each combination of the generated CV partitions, we perform a  
180 hyperparameter optimization on the current validation partition and evaluate the best-performing  
181 predictor on the current test partition. The experiments were also carried out using alternate regressors.  
182 See Sections N.1 and K for these results and all hyperparameter grids, respectively.

### 183 **3.4 Zero-shot predictors**

184 To investigate the efficacy of unsupervised learning on the curated datasets, we evaluate four zero-shot  
185 predictors. Using EVE, we evaluate the evidence lower bound (ELBO) by sampling and obtain a  
186 proxy for sequence fitness, analogous to the evolutionary index in [39]. Second and third proxies are  
187 obtained by evaluating the log-likelihood of a sequence conditioned on its structure using the inverse  
188 folding models ESM-IF1 [32] and ProteinMPNN [44]. The fourth zero-shot estimator is obtained  
189 by using Tranception [45] to evaluate the log-likelihood of each sequence. Details for the use of  
190 ProteinMPNN and Tranception can be found in Sections I and J in the supplementary materials.

Table 1: Summary of datasets and splits.

	$N_{tot}$	$N_A$	$N_B$	$N_C$	Split %ID	Target	Median %ID	Avg. length
<b>GH114</b>	55	20	18	17	0.55	Activity	0.46	268.8
<b>CM</b>	855	341	259	255	0.40	Activity	0.40	91.1
<b>PPAT</b>	615	182	234	199	0.55	Fitness	0.51	161.6

## 191 4 Datasets

192 The three curated datasets and the corresponding fitness landscapes are here motivated and described.  
 193 Despite the scarcity of available datasets described in Section 2, the curated datasets are representative  
 194 examples of wildtype discovery campaigns in terms of size and diversity. For additional curation  
 195 details on each dataset including specific thresholds for stratified splitting, see Section A.

### 196 4.1 GH114

197 **Motivation.** Accurately identifying enzymes with the highest activities towards a specific substrate  
 198 is of central importance during enzyme engineering. To achieve this, it is essential to ensure that assay  
 199 observations are directly comparable [46]. This includes maintaining identical experimental assay  
 200 conditions, including evaluating enzymes at the same concentrations and purity levels. However,  
 201 purifying enzymes requires significant work and resources, often resulting in assays composed of  
 202 fewer sequences, which are in turn of higher experimental quality.

203 **Landscape.** This dataset includes purified and concentration normalized natural glycoside hydro-  
 204 lase 114 (GH114) alpha-1,4-polygalactosaminidase enzymes and corresponding catalytic activity  
 205 values [47] which will act as the target of interest. GH114 enzymes degrade the exopolysaccharide  
 206 PEL, which provides structure and protection in some biofilms [48]. Having measurements of pu-  
 207 rified enzymes avoids issues with background effects from other enzymes in the recombinant host  
 208 background. We provide a curated version of the GH114 dataset which, to our knowledge, has not  
 209 been used in previous work for function prediction purposes.

### 210 4.2 CM

211 **Motivation.** Identification of enzymes with high catalytic activities is essential for enzyme engi-  
 212 neering campaigns. However, predicting the activity level of enzymes using physics-based methods  
 213 remains a great challenge [49]. Recent progress in high throughput screening allows the measurement  
 214 of enzyme activity of sequences with high diversity, but with low experimental cost.

215 **Landscape.** This dataset contains the catalytic activity of chorismate mutase (CM) homologous  
 216 proteins, as well as artificial sequences which follow the same pattern of variations (e.g., conservation  
 217 and co-evolution) [50]. The artificial sequences generated by Monte Carlo simulations at low and  
 218 medium temperatures match the empirical first-, second-, and higher-order statistics of the natural  
 219 homologs, while also exhibiting comparable catalytic levels when experimentally synthesized. These  
 220 sequence have therefore been included given the similarity in both sequence and fitness landscape.  
 221 See Section A.3 for further details. We perform an additional filtering of the dataset prior to the  
 222 splitting procedure by removing sequences with target values less than 0.42, corresponding to inactive  
 223 proteins [50]. This task thereby assumes that a preceding classification procedure has been carried  
 224 out. For completeness, we include benchmark results for the CM dataset when only the natural  
 225 homologs were used (see Section N.4) and classification results before the filtering step (see Section  
 226 M), which supports this last assumption.

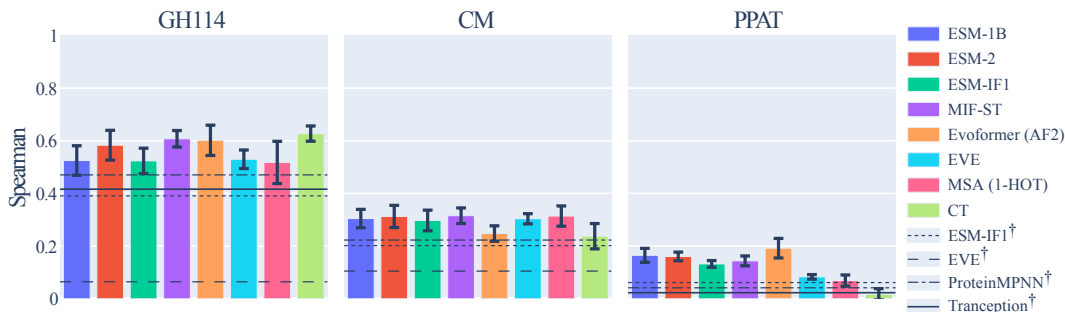


Figure 3: Average Spearman’s rank correlation (and standard error) between predictions and targets over test partitions. Higher is better. †: Zero-shot correlations.

### 227 4.3 PPAT

228 **Motivation.** PPAT (phosphopantetheine adenylyltransferase) is an essential enzyme that catalyzes  
 229 the second-to-last step in the CoA biosynthetic pathway. The target value for this prediction task  
 230 is the fitness score, which reflects the ability of PPAT homologs to complement a knockout E. coli  
 231 strain. The fitness of homologs can be affected by factors such as protein misfolding, mismatched  
 232 metabolic flux, or environmental mismatches etc. [51].

233 **Landscape.** This dataset contains fitness scores of 615 different PPAT homologs obtained by a  
 234 novel DNA synthesis/assembly technology, DropSynth, followed by a multiplexed functional assay  
 235 to measure how well each PPAT homolog can rescue a knockout phenotype [51].

### 236 4.4 Summary

237 A summary of the curated datasets including their total sizes, partition sizes, between-partition  
 238 sequence identity threshold (Split %ID), regression target, median pairwise sequence identity, and  
 239 average sequence length can be seen in Table 1. Additional curation details are found in the  
 240 supplementary materials (see Section A), while histograms of the regression targets for the three  
 241 datasets as well as partition histograms can be seen in Sections C and D, respectively.

242 The low median sequence identity observed in Table 1 highlights the diversity – and difficulty –  
 243 of wildtype datasets. Deep mutational scanning (DMS) datasets commonly used for variant effect  
 244 prediction on average have median sequence identities greater than 0.99. For comparison, Section F  
 245 shows both median, mean, and standard deviation for the curated datasets (see Table A1) contrasted  
 246 to 48 tasks from the ProteinGym benchmark (see Table A2). A visual example highlighting the  
 247 sequence diversity for a protein family compared to a DMS dataset can be seen in Figure 4 in [40].  
 248 The DMS data is localized to a small section of the protein family as it is composed of all single  
 249 mutations of a single wildtype. Similarly, a DMS of a protein belonging to the PPAT family would all  
 250 be positioned on a single branch of the phylogenetic tree in Figure 2.

## 251 5 Results

252 Spearman’s rank correlation between the predictions and targets over the three datasets can be seen in  
 253 Figure 3 and Table 2. The highest performing proteins are of interest making it the ranking and not  
 254 the absolute predictions that indicate the performance, despite the regressors being optimized using  
 255 the mean squared error. The RMSE can be found in Table A4 in the supplementary materials.

256 **GH114** The performance on the GH114 dataset highlights some of the peculiarities encountered  
 257 with small wildtype protein datasets. The collection of physicochemical descriptors (CT), which is  
 258 simpler than its competitors, achieves the highest score. Slightly below it are the structure-informed  
 259 MIF-ST and Evoformer representations, indicating a structural signal which is however not picked

Table 2: Benchmark results. Mean Spearman correlation and standard error using cross-validation. †: Zero-shot correlation on full datasets. Highest value and values within 1 SE are bold.

	<b>GH114</b>	<b>CM</b>	<b>PPAT</b>
ESM-1B	0.52 ± 0.06	<b>0.30</b> ± 0.03	<b>0.16</b> ± 0.03
ESM-2	<b>0.58</b> ± 0.06	<b>0.31</b> ± 0.04	<b>0.16</b> ± 0.02
ESM-IF1	0.52 ± 0.05	<b>0.30</b> ± 0.04	0.13 ± 0.01
MIF-ST	<b>0.61</b> ± 0.03	<b>0.32</b> ± 0.03	0.14 ± 0.02
Evoformer (AF2)	<b>0.60</b> ± 0.06	0.25 ± 0.03	<b>0.19</b> ± 0.04
EVE	0.53 ± 0.04	<b>0.30</b> ± 0.02	0.08 ± 0.01
MSA (1-HOT)	0.52 ± 0.08	<b>0.31</b> ± 0.04	0.07 ± 0.02
CT	<b>0.63</b> ± 0.03	0.24 ± 0.05	0.02 ± 0.02
ESM-IF1†	0.39	0.20	0.06
EVE†	0.06	0.11	−0.01
ProteinMPNN†	0.47	0.22	0.04
Tranception†	0.42	−0.05	0.02

up by the ESM-IF1 embeddings. While the CT representation achieves the highest mean value, several others are within one standard error, giving no clear advantage to neither complex nor simple models. While all supervised approaches beat the zero-shot predictors, ProteinMPNN, Tranception and ESM-IF1 likelihoods correlate well with the targets.

**CM** The second prediction task can be considered more challenging given the results, despite the comparatively large size of the CM dataset. While similar to the first task, an abundance of data is not sufficient to increase the downstream capabilities if it comes at the cost of potentially noisier measurements, as compared to the concentration normalized GH114 dataset. Most representations fall within one standard error of the top performer such that, once again, no representation paradigm has a clear advantage.

**PPAT** The most challenging task of the datasets, the results on the PPAT task show different behaviour. The evolutionary signal, i.e., the amount of information which can be learned from evolutionary homologs, is weak as indicated by the low correlations from the one-hot encoded MSA and from EVE (both in the supervised and zero-shot settings as per Table 2). Furthermore, the physicochemical descriptors fail to correlate – as do the remaining zero-shot predictors. The pLM and structure-based representations achieve the highest scores with the Evoformer embeddings coming out slightly ahead.

## 6 Ablation study

The dataset splitting procedure and benchmark tasks have been carefully constructed to ensure reliable estimates of model performance. In this section we show three ablation studies – one for each dataset – whereby different choices of task-structuring might lead to great over-estimations of performance. The results can be seen in Table 3 and in Figure A8. The  $\Delta$  columns in the table indicate differences to the benchmark results, where a positive/green value indicates *better* performance during ablation, i.e., over-estimation.

**Hold-out validation** For GH114, we perform hold-out validation by arbitrarily designating the three generated partitions as training, validation, and test sets and running the experiment only once. The correlations are significantly different to the benchmark results, with the ESM-2 correlation decreasing by 50 %. With no systematic pattern and decreased nuance given the lack of errorbars, it is easy to draw incorrect conclusions.

As the data-leakage between partitions has been controlled via the splitting procedure, the partitions are different from each other up to the sequence identity threshold. This implies that a model might



Table 3: Ablation results. Spearman correlation. \*: Hold-out validation, \*\*: Regression on both active and inactive proteins, \*\*\*: Repeated random splitting.  $\Delta$  shows difference to benchmark results. Highest value and values within 1 SE are bold.

	<b>GH114*</b>	$\Delta$	<b>CM**</b>	$\Delta$	<b>PPAT***</b>	$\Delta$
ESM-1B	0.36	-0.16	0.64 $\pm$ 0.01	+0.33	<b>0.23</b> $\pm$ 0.05	+0.07
ESM-2	0.39	-0.20	<b>0.66</b> $\pm$ 0.01	+0.35	0.18 $\pm$ 0.02	+0.02
ESM-IF1	0.46	-0.06	0.58 $\pm$ 0.01	+0.28	0.19 $\pm$ 0.01	+0.06
MIF-ST	0.62	+0.01	0.60 $\pm$ 0.02	+0.28	<b>0.24</b> $\pm$ 0.04	+0.09
Evoformer (AF2)	0.64	+0.04	0.57 $\pm$ 0.01	+0.33	<b>0.24</b> $\pm$ 0.02	+0.04
EVE	0.38 $\pm$ 0.04	-0.15	0.62 $\pm$ 0.00	+0.31	0.16 $\pm$ 0.01	+0.08
MSA (1-HOT)	0.50	-0.02	0.61 $\pm$ 0.01	+0.30	0.06 $\pm$ 0.06	-0.01
CT	<b>0.65</b>	+0.03	0.52 $\pm$ 0.01	+0.28	0.13 $\pm$ 0.01	+0.11

291 perform well on, e.g., only a subset of the partitions. Choosing which partitions to use for training,  
 292 validation, and testing is (in this case) arbitrary and can thereby lead to misleading results. To  
 293 avoid this pitfall, cross-validation is needed such that the average predictive performance on all  
 294 combinations of partitions can be estimated. An analogue ablation study for the CM and PPAT  
 295 datasets can be found in Section L.1, where similar conclusions can be drawn.

296 **Disregarding distinct target modalities** For the CM dataset, we only included the active sequences  
 297 in the benchmark. To demonstrate why, we have included the results of performing regression on  
 298 both active *and* inactive sequences in the center of Table 3. These results are greatly overinflated  
 299 compared to the benchmark results, with some representations more than doubling the correlation  
 300 scores.

301 Regression performed on a dataset with a distinctly bimodal target distribution (such as the full  
 302 CM dataset, see Figure C in the supplementary materials) can inflate the results significantly. The  
 303 regressor is able to distinguish between the two target modalities, i.e., between the inactive cluster  
 304 around 0 and the active cluster around 1, driving the ranking correlation to overly-optimistic values.  
 305 The caveat to this preprocessing step is that it requires knowing the whether the proteins are active  
 306 or not a priori, which assumes that a preceding classification-screening has been performed. The  
 307 classification results of such a process can be seen in Section M.

308 **Random partitioning for cross-validation** To illustrate why random splitting of wildtype protein  
 309 datasets is ill-advised, we applied repeated random splitting to the PPAT dataset. This was done by  
 310 randomly assigning sequences to training, validation, and testing partitions without any consideration  
 311 of sequence similarity. Given the randomized partitions, the predictive performance using the selected  
 312 representations was evaluated using cross validation. This was repeated a total of three times with  
 313 different seeds. While the results look similar to the benchmark results, we do see an increase in  
 314 performance across the board.

315 With random sampling, we risk placing very similar sequences in separate partitions, thereby allowing  
 316 extensive data-leakage, where we are essentially testing on training/validation data, thus overestimating  
 317 the predictive performance [52]. The results for this ablation study carried out on the GH114 and  
 318 CM datasets can be found in Section L.2, where we can once again draw similar conclusions.

## 319 7 Discussion

320 The choice of representation greatly affects the downstream predictive capabilities, with no consistent,  
 321 clear edge given by any of the three representation paradigms. For CM, a one-hot encoded MSA acts as  
 322 an impressive baseline proving difficult to convincingly beat. For GH114, physicochemical descriptors  
 323 are sufficient to achieve top performance, while the PPAT dataset benefits from the complex, structure-  
 324 informed Evoformer embeddings. While the specific top-scoring representation fluctuates, the  
 325 ESM-2 embeddings are consistently within one standard error and can thus be considered a relatively

326 consistent baseline for future experiments, where others occasionally underperform. For the three  
327 tasks, we see supervised learning outperforming zero-shot predictions, while the inverse-folding  
328 estimators however offer decent zero-shot approaches for two of three tasks.

329 Despite similar overall patterns, some results stand out, e.g., the comparatively high performance  
330 on the GH114 dataset and the low performance on the PPAT dataset. Variations in experimental  
331 conditions and techniques can introduce different levels of noise. The CM and PPAT datasets are  
332 derived from tests on supernatants with complex backgrounds with potential side-activities from  
333 impurities, whereas the GH114 dataset uses purified samples with less expected noise. This can be a  
334 potential reason for the comparatively high performance of the latter. As for the low performance on  
335 the PPAT dataset, the reason might lie in the target values: the GH114 and CM datasets both measure  
336 enzymatic activities while the PPAT dataset measures fitness. The overall performance disparities  
337 suggest that enzyme activities, rather than a more complex and assay-specific fitness value, are easier  
338 to model given the available protein representation paradigms. The stark contrast in performance  
339 between the concentration normalized GH114 dataset and both the CM and PPAT datasets indicates  
340 that higher quality datasets are of central importance to learn accurate fitness landscapes – more so  
341 than the number of labelled sequences.

## 342 **8 Conclusion**

343 In this work we have presented a novel benchmark which investigates an unexplored domain of  
344 machine learning-driven protein engineering: the navigation of global fitness landscapes for single  
345 protein families. Wildtype exploration can be viewed as a predominantly explorative phase of  
346 protein optimization, which precedes the exploitation phase comprised of the subsequent protein  
347 engineering. Often, limited resources are allocated to wildtype exploration since it is inherently  
348 costly and considered wasteful as it tends to produce many poor candidates. This is unlikely to  
349 change unless we find ways to improve our wildtype search strategy, which will require better  
350 predictions. We therefore consider the limited dataset sizes as an inherent condition and limitation in  
351 this domain. This makes the collection and curation of relevant labelled datasets challenging and also  
352 necessitates the design of careful learning schemes and model evaluation to ensure reliable estimates  
353 of generalizability while avoiding inadvertently overestimating the results. We anticipate that the  
354 creation of this new set of comprehensive family-wide datasets will facilitate and improve future  
355 model development and applicability in this domain.

356 Given the limited dataset sizes, our focus has been on transfer learning and zero-shot prediction. Our  
357 results show that the supervised approaches outperform the zero-shot approaches, but that no one  
358 representation or representation paradigm consistently outperforms the others. This could suggest  
359 that the employed representations are not sufficiently informative. A key limitation for a number  
360 of the included representations is that we obtained protein-level representations as averages over  
361 the protein length to arrive at fixed-length embeddings, which is known to be suboptimal [28]. We  
362 encourage the community to experiment with novel aggregation strategies and new representation  
363 designs to improve performance on our benchmark. It is also conceivable that general-purpose  
364 protein representation models might not by themselves be sufficient to convincingly improve on the  
365 proposed tasks. One can imagine that further improvements can be obtained using pretrained models  
366 fine-tuned on a protein family of interest – or by developing weakly-supervised representation models  
367 incorporating relevant properties that correlate with the function of interest (e.g., thermostability).

368 Although the performance of current baselines on some of our test-cases is fairly low in absolute  
369 terms, even low correlations can provide useful guidance on selecting wildtype protein starting points  
370 and can have measurable real-world impacts. Any further improvements will enhance the importance  
371 of wildtype exploration relative to the subsequent local optimization step. In silico screenings of  
372 potential wildtype candidates can be scaled efficiently compared to expensive, time-consuming in  
373 vitro assays, significantly reducing the early costs of future protein engineering campaigns. We hope  
374 that FLOP will pave the way for these developments.

## 375 **Acknowledgments and Disclosure of Funding**

376 This work was funded in part by Innovation Fund Denmark (1044-00158A), the Novo Nordisk  
377 Foundation through the MLSS Center (Basic Machine Learning Research in Life Science,  
378 NNF200C0062606), the Pioneer Centre for AI (DRNF grant number P1), and the Danish Data  
379 Science Academy, which is funded by the Novo Nordisk Foundation (NNF21SA0069429) and  
380 VILLUM FONDEN (40516).

## 381 **References**

- 382 [1] B. J. Wittmann, Y. Yue, and F. H. Arnold, “Informed training set design enables efficient machine  
383 learning-assisted directed protein evolution,” *Cell systems*, vol. 12, no. 11, pp. 1026–1045, 2021.
- 384 [2] T. Davids, M. Schmidt, D. Böttcher, and U. T. Bornscheuer, “Strategies for the discovery and  
385 engineering of enzymes for biocatalysis,” *Current Opinion in Chemical Biology*, vol. 17, no. 2,  
386 pp. 215–220, 2013.
- 387 [3] K. K. Yang, Z. Wu, and F. H. Arnold, “Machine-learning-guided directed evolution for protein  
388 engineering,” *Nature methods*, vol. 16, no. 8, pp. 687–694, 2019.
- 389 [4] P. E. Bourne, “Casp and casp experiments and their findings,” *Methods of Biochemical*  
390 *Analysis*, vol. 44, pp. 501–508, 2003.
- 391 [5] J. Moulton, “A decade of casp: progress, bottlenecks and prognosis in protein structure prediction,”  
392 *Current opinion in structural biology*, vol. 15, no. 3, pp. 285–289, 2005.
- 393 [6] A. Kryzhtafovich, T. Schwede, M. Topf, K. Fidelis, and J. Moulton, “Critical assessment of  
394 methods of protein structure prediction (casp)—round xiv,” *Proteins: Structure, Function, and*  
395 *Bioinformatics*, vol. 89, no. 12, pp. 1607–1617, 2021.
- 396 [7] R. Rao, N. Bhattacharya, N. Thomas, Y. Duan, X. Chen, J. Canny, P. Abbeel, and Y. S. Song,  
397 “Evaluating Protein Transfer Learning with TAPE,” *arXiv:1906.08230 [cs, q-bio, stat]*, Jun.  
398 2019, arXiv: 1906.08230. [Online]. Available: <http://arxiv.org/abs/1906.08230>
- 399 [8] G. J. Rocklin, T. M. Chidyausiku, I. Goreshnik, A. Ford, S. Houlston, A. Lemak, L. Carter,  
400 R. Ravichandran, V. K. Mulligan, A. Chevalier, C. H. Arrowsmith, and D. Baker, “Global  
401 analysis of protein folding using massively parallel design, synthesis, and testing,” *Science (New*  
402 *York, N.Y.)*, vol. 357, no. 6347, pp. 168–175, Jul. 2017.
- 403 [9] K. S. Sarkisyan, D. A. Bolotin, M. V. Meer, D. R. Usmanova, A. S. Mishin, G. V. Sharonov,  
404 D. N. Ivankov, N. G. Bozhanova, M. S. Baranov, O. Soylemez *et al.*, “Local fitness landscape  
405 of the green fluorescent protein,” *Nature*, vol. 533, no. 7603, pp. 397–401, 2016.
- 406 [10] M. Xu, Z. Zhang, J. Lu, Z. Zhu, Y. Zhang, C. Ma, R. Liu, and J. Tang, “PEER: A Comprehensive  
407 and Multi-Task Benchmark for Protein Sequence Understanding,” Jun. 2022, number:  
408 arXiv:2206.02096 arXiv:2206.02096 [cs]. [Online]. Available: <http://arxiv.org/abs/2206.02096>
- 409 [11] V. E. Gray, R. J. Hause, J. Luebeck, J. Shendure, and D. M. Fowler, “Quantitative Missense  
410 Variant Effect Prediction Using Large-Scale Mutagenesis Data,” *Cell Systems*, vol. 6, no. 1, pp.  
411 116–124.e3, Jan. 2018.
- 412 [12] P. Notin, M. Dias, J. Frazer, J. Marchena-Hurtado, A. Gomez, D. S. Marks, and Y. Gal,  
413 “Tranception: protein fitness prediction with autoregressive transformers and inference-time  
414 retrieval,” May 2022, number: arXiv:2205.13760 arXiv:2205.13760 [cs]. [Online]. Available:  
415 <http://arxiv.org/abs/2205.13760>
- 416 [13] R. J. L. Townshend, M. Vögele, P. Suriana, A. Derry, A. Powers, Y. Laloudakis, S. Balachandar,  
417 B. Jing, B. Anderson, S. Eismann, R. Kondor, R. B. Altman, and R. O. Dror, “ATOM3D: Tasks  
418 On Molecules in Three Dimensions,” *arXiv:2012.04035 [physics, q-bio]*, Jan. 2022, arXiv:  
419 2012.04035. [Online]. Available: <http://arxiv.org/abs/2012.04035>

- 420 [14] I. H. Moal and J. Fernández-Recio, “SKEMPI: a Structural Kinetic and Energetic database of  
421 Mutant Protein Interactions and its use in empirical models,” *Bioinformatics*, vol. 28, no. 20, pp.  
422 2600–2607, Oct. 2012. [Online]. Available: <https://doi.org/10.1093/bioinformatics/bts489>
- 423 [15] J. Jankauskaite, B. Jiménez-García, J. Dapkunas, J. Fernández-Recio, and I. H. Moal, “SKEMPI  
424 2.0: an updated benchmark of changes in protein-protein binding energy, kinetics and thermo-  
425 dynamics upon mutation,” *Bioinformatics (Oxford, England)*, vol. 35, no. 3, pp. 462–469, Feb.  
426 2019.
- 427 [16] C. Dallago, J. Mou, K. E. Johnston, B. J. Wittmann, N. Bhattacharya, S. Goldman,  
428 A. Madani, and K. K. Yang, “FLIP: Benchmark tasks in fitness landscape inference for  
429 proteins,” Jan. 2022, pages: 2021.11.09.467890 Section: New Results. [Online]. Available:  
430 <https://www.biorxiv.org/content/10.1101/2021.11.09.467890v2>
- 431 [17] A. Jarzab, N. Kurzawa, T. Hopf, M. Moerch, J. Zecha, N. Leijten, Y. Bian, E. Musiol,  
432 M. Maschberger, G. Stoehr, I. Becher, C. Daly, P. Samaras, J. Mergner, B. Spanier, A. Angelov,  
433 T. Werner, M. Bantscheff, M. Wilhelm, M. Klingenspor, S. Lemeer, W. Liebl, H. Hahne,  
434 M. M. Savitski, and B. Kuster, “Meltome atlas—thermal proteome stability across the tree of  
435 life,” *Nature Methods*, vol. 17, no. 5, pp. 495–503, May 2020, number: 5 Publisher: Nature  
436 Publishing Group. [Online]. Available: <https://www.nature.com/articles/s41592-020-0801-4>
- 437 [18] N. C. Wu, L. Dai, C. A. Olson, J. O. Lloyd-Smith, and R. Sun, “Adaptation in protein fitness  
438 landscapes is facilitated by indirect paths,” *Elife*, vol. 5, p. e16965, 2016.
- 439 [19] D. H. Bryant, A. Bashir, S. Sinai, N. K. Jain, P. J. Ogden, P. F. Riley, G. M. Church, L. J.  
440 Colwell, and E. D. Kelsic, “Deep diversification of an aav capsid protein by machine learning,”  
441 *Nature Biotechnology*, vol. 39, no. 6, pp. 691–696, 2021.
- 442 [20] M. H. Gíslason, H. Nielsen, J. J. Almagro Armenteros, and A. R. Johansen, “Prediction of  
443 gpi-anchored proteins with pointer neural networks,” *Current Research in Biotechnology*,  
444 vol. 3, pp. 6–13, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2590262821000010>
- 446 [21] F. Teufel, M. H. Gíslason, J. J. A. Armenteros, A. R. Johansen, O. Winther, and  
447 H. Nielsen, “GraphPart: Homology partitioning for biological sequence analysis,”  
448 Apr. 2023, pages: 2023.04.14.536886 Section: New Results. [Online]. Available:  
449 <https://www.biorxiv.org/content/10.1101/2023.04.14.536886v1>
- 450 [22] M. Steinegger and J. Söding, “MMseqs2 enables sensitive protein sequence searching  
451 for the analysis of massive data sets,” *Nature Biotechnology*, vol. 35, no. 11, pp.  
452 1026–1028, Nov. 2017, number: 11 Publisher: Nature Publishing Group. [Online]. Available:  
453 <https://www.nature.com/articles/nbt.3988>
- 454 [23] K. K. Yang, A. X. Lu, and N. Fusi, “Convolutions are competitive with transformers for protein  
455 sequence pretraining,” May 2022, pages: 2022.05.19.492714 Section: New Results. [Online].  
456 Available: <https://www.biorxiv.org/content/10.1101/2022.05.19.492714v2>
- 457 [24] A. Rives, J. Meier, T. Sercu, S. Goyal, Z. Lin, J. Liu, D. Guo, M. Ott, C. L. Zitnick, J. Ma, and  
458 R. Fergus, “Biological structure and function emerge from scaling unsupervised learning to 250  
459 million protein sequences,” *bioRxiv, Tech. Rep.*, Dec. 2020, section: New Results Type: article.  
460 [Online]. Available: <https://www.biorxiv.org/content/10.1101/622803v4>
- 461 [25] J. Meier, R. Rao, R. Verkuil, J. Liu, T. Sercu, and A. Rives, “Language models  
462 enable zero-shot prediction of the effects of mutations on protein function,” *bioRxiv*,  
463 *Tech. Rep.*, Jul. 2021, section: New Results Type: article. [Online]. Available:  
464 <https://www.biorxiv.org/content/10.1101/2021.07.09.450648v1>

- 465 [26] R. M. Rao, J. Liu, R. Verkuil, J. Meier, J. Canny, P. Abbeel, T. Sercu, and A. Rives,  
466 “MSA Transformer,” in *Proceedings of the 38th International Conference on Machine*  
467 *Learning*. PMLR, Jul. 2021, pp. 8844–8856, iSSN: 2640-3498. [Online]. Available:  
468 <https://proceedings.mlr.press/v139/rao21a.html>
- 469 [27] Z. Lin, H. Akin, R. Rao, B. Hie, Z. Zhu, W. Lu, A. d. Santos Costa, M. Fazel-Zarandi,  
470 T. Sercu, S. Candido, and A. Rives, “Language models of protein sequences at the scale of  
471 evolution enable accurate structure prediction,” *Synthetic Biology*, preprint, Jul. 2022. [Online].  
472 Available: <http://biorxiv.org/lookup/doi/10.1101/2022.07.20.500902>
- 473 [28] N. S. Detlefsen, S. Hauberg, and W. Boomsma, “Learning meaningful representations  
474 of protein sequences,” *Nature Communications*, vol. 13, no. 1, p. 1914, Apr.  
475 2022, number: 1 Publisher: Nature Publishing Group. [Online]. Available: <https://www.nature.com/articles/s41467-022-29443-w>
- 477 [29] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool,  
478 R. Bates, A. Žídek, A. Potapenko, A. Bridgland, C. Meyer, S. A. A. Kohl, A. J. Ballard,  
479 A. Cowie, B. Romera-Paredes, S. Nikolov, R. Jain, J. Adler, T. Back, S. Petersen, D. Reiman,  
480 E. Clancy, M. Zielinski, M. Steinegger, M. Pacholska, T. Berghammer, S. Bodenstein,  
481 D. Silver, O. Vinyals, A. W. Senior, K. Kavukcuoglu, P. Kohli, and D. Hassabis, “Highly  
482 accurate protein structure prediction with AlphaFold,” *Nature*, vol. 596, no. 7873, pp.  
483 583–589, Aug. 2021, number: 7873 Publisher: Nature Publishing Group. [Online]. Available:  
484 <https://www.nature.com/articles/s41586-021-03819-2>
- 485 [30] M. Mirdita, K. Schütze, Y. Moriwaki, L. Heo, S. Ovchinnikov, and M. Steinegger,  
486 “ColabFold: making protein folding accessible to all,” *Nature Methods*, vol. 19, no. 6, pp.  
487 679–682, Jun. 2022, number: 6 Publisher: Nature Publishing Group. [Online]. Available:  
488 <https://www.nature.com/articles/s41592-022-01488-1>
- 489 [31] M. Hu, F. Yuan, K. K. Yang, F. Ju, J. Su, H. Wang, F. Yang, and Q. Ding, “Exploring  
490 evolution-based & -free protein language models as protein function predictors,” Jun.  
491 2022, number: arXiv:2206.06583 arXiv:2206.06583 [cs, q-bio]. [Online]. Available:  
492 <http://arxiv.org/abs/2206.06583>
- 493 [32] C. Hsu, R. Verkuil, J. Liu, Z. Lin, B. Hie, T. Sercu, A. Lerer, and A. Rives, “Learning  
494 inverse folding from millions of predicted structures,” in *Proceedings of the 39th International*  
495 *Conference on Machine Learning*. PMLR, Jun. 2022, pp. 8946–8970, iSSN: 2640-3498.  
496 [Online]. Available: <https://proceedings.mlr.press/v162/hsu22a.html>
- 497 [33] K. K. Yang, H. Yeh, and N. Zanichelli, “Masked inverse folding with sequence transfer for  
498 protein representation learning,” Mar. 2023, pages: 2022.05.25.493516 Section: New Results.  
499 [Online]. Available: <https://www.biorxiv.org/content/10.1101/2022.05.25.493516v3>
- 500 [34] F. Morcos, A. Pagnani, B. Lunt, A. Bertolino, D. S. Marks, C. Sander, R. Zecchina, J. N.  
501 Onuchic, T. Hwa, and M. Weigt, “Direct-coupling analysis of residue coevolution captures  
502 native contacts across many protein families,” *Proceedings of the National Academy of Sciences*,  
503 vol. 108, no. 49, pp. E1293–E1301, 2011.
- 504 [35] D. S. Marks, L. J. Colwell, R. Sheridan, T. A. Hopf, A. Pagnani, R. Zecchina, and C. Sander,  
505 “Protein 3d structure computed from evolutionary sequence variation,” *PloS one*, vol. 6, no. 12,  
506 p. e28766, 2011.
- 507 [36] P. Tian and R. B. Best, “How many protein sequences fold to a given structure? a coevolutionary  
508 analysis,” *Biophysical journal*, vol. 113, no. 8, pp. 1719–1730, 2017.
- 509 [37] The UniProt Consortium, “UniProt: the Universal Protein Knowledgebase in 2023,”  
510 *Nucleic Acids Research*, vol. 51, no. D1, pp. D523–D531, Jan. 2023. [Online]. Available:  
511 <https://doi.org/10.1093/nar/gkac1052>

- 512 [38] M. Blum, H.-Y. Chang, S. Chuguransky, T. Grego, S. Kandasamy, A. Mitchell, G. Nuka,  
513 T. Paysan-Lafosse, M. Qureshi, S. Raj, L. Richardson, G. A. Salazar, L. Williams, P. Bork,  
514 A. Bridge, J. Gough, D. H. Haft, I. Letunic, A. Marchler-Bauer, H. Mi, D. A. Natale, M. Necci,  
515 C. A. Orengo, A. P. Pandurangan, C. Rivoire, C. J. A. Sigrist, I. Sillitoe, N. Thanki, P. D.  
516 Thomas, S. C. E. Tosatto, C. H. Wu, A. Bateman, and R. D. Finn, “The interpro protein families  
517 and domains database: 20 years on,” *Nucleic acids research*, vol. 49, no. D1, p. D344–D354,  
518 January 2021. [Online]. Available: <https://europepmc.org/articles/PMC7778928>
- 519 [39] J. Frazer, P. Notin, M. Dias, A. Gomez, J. K. Min, K. Brock, Y. Gal, and D. S. Marks, “Disease  
520 variant prediction with deep generative models of evolutionary data,” *Nature*, vol. 599, no.  
521 7883, pp. 91–95, Nov. 2021, number: 7883 Publisher: Nature Publishing Group. [Online].  
522 Available: <https://www.nature.com/articles/s41586-021-04043-8>
- 523 [40] A. J. Riesselman, J. B. Ingraham, and D. S. Marks, “Deep generative models of  
524 genetic variation capture the effects of mutations,” *Nature Methods*, vol. 15, no. 10, pp.  
525 816–822, Oct. 2018, number: 10 Publisher: Nature Publishing Group. [Online]. Available:  
526 <https://www.nature.com/articles/s41592-018-0138-4>
- 527 [41] T. A. Hopf, A. G. Green, B. Schubert, S. Mersmann, C. P. I. Schärfe, J. B. Ingraham, A. Toth-  
528 Petroczy, K. Brock, A. J. Riesselman, P. Palmedo, C. Kang, R. Sheridan, E. J. Draizen, C. Dal-  
529 lago, C. Sander, and D. S. Marks, “The EVcouplings Python framework for coevolutionary  
530 sequence analysis,” *Bioinformatics (Oxford, England)*, vol. 35, no. 9, pp. 1582–1584, May  
531 2019.
- 532 [42] D. Hesslow, N. Zanichelli, P. Notin, I. Poli, and D. Marks, “RITA: a Study on Scaling  
533 Up Generative Protein Sequence Models,” arXiv, Tech. Rep. arXiv:2205.05789, May 2022,  
534 arXiv:2205.05789 [cs, q-bio] type: article. [Online]. Available: <http://arxiv.org/abs/2205.05789>
- 535 [43] J. Dong, Z.-J. Yao, L. Zhang, F. Luo, Q. Lin, A.-P. Lu, A. F. Chen, and D.-S. Cao, “PyBioMed:  
536 a python library for various molecular representations of chemicals, proteins and DNAs and  
537 their interactions,” *Journal of Cheminformatics*, vol. 10, no. 1, p. 16, Mar. 2018.
- 538 [44] J. Dauparas, I. Anishchenko, N. Bennett, H. Bai, R. J. Ragotte, L. F. Milles, B. I. M.  
539 Wicky, A. Courbet, R. J. d. Haas, N. Bethel, P. J. Y. Leung, T. F. Huddy, S. Pellock,  
540 D. Tischer, F. Chan, B. Koepnick, H. Nguyen, A. Kang, B. Sankaran, A. K. Bera,  
541 N. P. King, and D. Baker, “Robust deep learning based protein sequence design using  
542 ProteinMPNN,” Jun. 2022, pages: 2022.06.03.494563 Section: New Results. [Online].  
543 Available: <https://www.biorxiv.org/content/10.1101/2022.06.03.494563v1>
- 544 [45] P. Notin, M. Dias, J. Frazer, J. M. Hurtado, A. N. Gomez, D. Marks, and Y. Gal, “Tranception:  
545 protein fitness prediction with autoregressive transformers and inference-time retrieval,” in  
546 *International Conference on Machine Learning*. PMLR, 2022, pp. 16990–17017.
- 547 [46] R. Lonsdale, J. N. Harvey, and A. J. Mulholland, “A practical guide to modelling  
548 enzyme-catalysed reactions,” *Chemical Society Reviews*, vol. 41, no. 8, pp. 3025–3038, Apr.  
549 2012. [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3371381/>
- 550 [47] M. Li, J. Salomon, D. R. Segura, M. A. Stringer, R. M. Vejborg, D. M. K. Klitgaard, D. Nissen,  
551 W. Peng, and T. Sun, “Polypeptides,” Patent WO/2019/228448, Dec., 2019. [Online]. Available:  
552 <https://patentscope.wipo.int/search/en/detail.jsf?docId=WO2019228448>
- 553 [48] K. M. Colvin, V. D. Gordon, K. Murakami, B. R. Borlee, D. J. Wozniak, G. C. L.  
554 Wong, and M. R. Parsek, “The Pel Polysaccharide Can Serve a Structural and Protective  
555 Role in the Biofilm Matrix of *Pseudomonas aeruginosa*,” *PLOS Pathogens*, vol. 7,  
556 no. 1, p. e1001264, Jan. 2011, publisher: Public Library of Science. [Online]. Available:  
557 <https://journals.plos.org/plospathogens/article?id=10.1371/journal.ppat.1001264>

- 558 [49] M. K. Tiwari, R. Singh, R. K. Singh, I.-W. Kim, and J.-K. Lee, “Computational approaches for  
559 rational design of proteins with novel functionalities,” *Computational and structural biotechnol-  
560 ogy journal*, vol. 2, no. 3, p. e201204002, 2012.
- 561 [50] W. P. Russ, M. Figliuzzi, C. Stocker, P. Barrat-Charlaix, M. Socolich, P. Kast, D. Hilvert,  
562 R. Monasson, S. Cocco, M. Weigt, and R. Ranganathan, “An evolution-based model for  
563 designing chorismate mutase enzymes,” *Science*, vol. 369, no. 6502, pp. 440–445, Jul. 2020,  
564 publisher: American Association for the Advancement of Science. [Online]. Available:  
565 <https://www.science.org/doi/full/10.1126/science.aba3304>
- 566 [51] C. Plesa, A. M. Sidore, N. B. Lubock, D. Zhang, and S. Kosuri, “Multiplexed gene synthesis  
567 in emulsions for exploring protein functional landscapes,” *Science*, vol. 359, no. 6373, pp.  
568 343–347, Jan. 2018, publisher: American Association for the Advancement of Science.  
569 [Online]. Available: <https://www.science.org/doi/10.1126/science.aao5167>
- 570 [52] P. Bork and E. V. Koonin, “Predicting functions from protein sequences—where are the bottle-  
571 necks?” *Nature Genetics*, vol. 18, no. 4, pp. 313–318, Apr. 1998.

572 **Checklist**

- 573 1. For all authors...
- 574 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
575 contributions and scope? [Yes]
- 576 (b) Did you describe the limitations of your work? [Yes]
- 577 (c) Did you discuss any potential negative societal impacts of your work? [No] We do  
578 not believe that there are any potential negative societal impacts of this work, since it  
579 concerns modelling naturally occurring wildtype proteins.
- 580 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
581 them? [Yes]
- 582 2. If you are including theoretical results...
- 583 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 584 (b) Did you include complete proofs of all theoretical results? [N/A]
- 585 3. If you ran experiments (e.g. for benchmarks)...
- 586 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
587 mental results (either in the supplemental material or as a URL)? [Yes] See supplement-  
588 ary materials.
- 589 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
590 were chosen)? [Yes] See both Section 3.1 for data splits and supplementary materials  
591 Section K for hyperparameter grids.
- 592 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
593 ments multiple times)? [Yes]
- 594 (d) Did you include the total amount of compute and the type of resources used (e.g.,  
595 type of GPUs, internal cluster, or cloud provider)? [Yes] See supplementary materials,  
596 Section B.1.
- 597 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 598 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 599 (b) Did you mention the license of the assets? [No] No specific license agreements were  
600 found. Consent was however given from authors of the datasets to include them in the  
601 benchmark, and we specify in the supplementary materials Section A that references to  
602 tasks in this benchmark should include references to the original data sources.
- 603 (c) Did you include any new assets either in the supplemental material or as a URL? [No]  
604 All datasets and methods are publicly available prior to this work. We do however  
605 provide thorough descriptions of how to access and use the original data as well as our  
606 processed/curated versions of it.
- 607 (d) Did you discuss whether and how consent was obtained from people whose data you're  
608 using/curating? [Yes] Consent was given explicitly for all datasets.
- 609 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
610 information or offensive content? [N/A]
- 611 5. If you used crowdsourcing or conducted research with human subjects...
- 612 (a) Did you include the full text of instructions given to participants and screenshots, if  
613 applicable? [N/A]
- 614 (b) Did you describe any potential participant risks, with links to Institutional Review  
615 Board (IRB) approvals, if applicable? [N/A]
- 616 (c) Did you include the estimated hourly wage paid to participants and the total amount  
617 spent on participant compensation? [N/A]