UNLEARNING VIRUS KNOWLEDGE TOWARD SAFE AND RESPONSIBLE MUTATION EFFECT PREDICTIONS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

024

025 026

027

037

Paper under double-blind review

ABSTRACT

Pre-trained deep protein models have become essential tools in fields such as biomedical research, enzyme engineering, and therapeutics due to their ability to predict and optimize protein properties effectively. However, the diverse and broad training data used to enhance the generalizability of these models may also inadvertently introduce ethical risks and pose biosafety concerns, such as the enhancement of harmful viral properties like transmissibility or drug resistance. To address this issue, we introduce a novel approach using knowledge unlearning to selectively remove virus-related knowledge while retaining other useful capabilities. We propose a learning scheme, PROEDIT, for editing a pre-trained protein language model toward safe and responsible mutation effect prediction. Extensive validation on open benchmarks demonstrates that PROEDIT significantly reduces the model's ability to enhance the properties of virus mutants without compromising its performance on non-virus proteins. As the first thorough exploration of safety issues in deep learning solutions for protein engineering, this study provides a foundational step toward ethical and responsible AI in biology.

1 INTRODUCTION

Pre-trained deep protein models are playing an increasingly important role in biological research (Narayanan et al., 2021; Pucci et al., 2022). By learning from massive amounts of existing protein data, these models uncover hidden relationships between protein sequences, structures, functions, and dynamics. Remarkable successes have been witnessed in diverse applications, such as enzyme design (Madani et al., 2023; Zhou et al., 2024).

Similar to natural language processing, pre-trained protein models often require training on billions of sequences with large-scale models to enhance expressivity and generalizability, achieving top performance across downstream tasks (Laine et al., 2019; Notin et al., 2022b; Lin et al., 2023b). This framework has been widely applied in solving problems in molecule design, where labels are usually scarce, expensive, or nonexistent. For instance, in enzyme engineering, *mutation effect prediction* (Notin



Figure 1: Positive Relationship of A model's overall performance vs Virusrelated performance on mutation effect prediction. Data source: https:// proteingym.org/benchmarks

et al., 2024) uses pre-trained models to score and rank the fitness of mutants relative to arbitrary
wild-type proteins. Deep learning models guide proteins to modify toward enhanced functionalities such as activity, stability, and yield. Compared to previous rational design or simulation-based
methods, they significantly improve mutation design success rates and reduce experimental costs
by recommending better mutation strategies, while not relying on specific biological knowledge or
experimental data (Lu et al., 2022; Li et al., 2023; Zhou et al., 2024a).

As is well known, a model's output is directly influenced by its training data. To improve expressivity and generalization, pre-trained models typically incorporate a large and diverse dataset to learn the parameters of the model. However, some of this knowledge may inevitably contain factual er-

069

071 072

054 /irus $\mathbb{D}^{\mathrm{fgt}}$ Mode Unlearn Scope 056 virus know Forget Virus $\max - \mathbb{E}(\mathbf{x}_{\text{virus}})$ Knowledge Language Model uage Model Language 058 $\mathbb{D}^{\mathrm{sim}}$ Retention Scope /irus-like orotein knowledge 060 -Virus $\max \mathbb{E}(\mathbf{x}_{non-virus})$ 061 Protein t 062 Corruption Scope Distinguish us vs non-Virus Knowledge 063 $\mathbb{D}^{\mathrm{norm}}$ 064 max KL(PLM||PLM° protein knowledge Edited 065 066 non-virus Model Unlearning 067

Figure 2: Illustration of the proposed PROEDIT. A pre-trained PLM is updated using three datasets to selectively forget virus-related knowledge while retaining non-virus information. To further enhance PROEDIT's ability to distinguish between virus and non-virus data, an additional dataset of 070 virus-like proteins is specifically prepared to guide the optimization.

073 rors, biases, or even harmful information. These misleading elements can severely undermine the 074 reliability and ethical integrity of the content generated by models pre-trained on such data. This 075 issue is prevalent across many domains, such as natural language processing (Wang et al., 2024b) 076 and computer vision (Golatkar et al., 2021). 077

The same concern has been raised in mutation effect prediction tasks. Figure 1 analyzes the 078 performance of existing models on ProteinGym leaderboard (https://proteingym.org/ 079 benchmarks) for mutation effect prediction. It is evident that these powerful models show a strong positive correlation between reliability in modifying arbitrary enzymes and viruses. Designing tools 081 capable of enhancing viral properties (e.g., transmission, immune evasion, and drug resistance) and 082 offering them to the public poses significant biosafety and ethical risks, such as disrupting ecological 083 balance, triggering severe pandemics, and fostering biological weapons. 084

Therefore, it is urgent and important to develop corresponding techniques to edit protein models and 085 allow the final models to retain their ability to effectively improve other enzymes while significantly 086 reducing their capacity to enhance viruses, thereby mitigating ethical risks and enhancing the safety 087 of research. While this issue has been preliminarily discussed in recent studies (Truong Jr & Bepler, 880 2023; Tan et al., 2023; Ouyang-Zhang et al., 2024; Liu et al., 2024), to the best of our knowledge, 089 no solution has been developed to edit the pre-trained model and address this problem. 090

To this end, we employ the knowledge unlearning technique (Sinitsin et al., 2019; Wang et al., 2023) 091 and propose PROEDIT, a learning scheme for safe and responsible protein language models (PLMs) 092 for mutation effect prediction. We distinguish three types of data from the UniRef database: "virus", 093 "non-virus", and "virus-like non-virus", and construct corresponding learning objectives. A pre-094 trained PLM is guided to retain its understanding of non-virus data within the retention scope while 095 forgetting virus-related information within the unlearning scope (Figure 2). Notably, we introduce 096 an additional corruption scope to ensure that the unlearned model retains the ability to understand virus-like non-virus data. Empirically, we validate that PROEDIT significantly reduces prediction 098 performance on virus mutants across various virus assays while maintaining strong performance on non-virus mutants. In contrast, existing models either improve or degrade performance on both virus 099 and non-virus assays simultaneously. 100

101 In summary, this work addresses mutation effect prediction-a core challenge in protein engineering-102 and presents the first detailed discussion of safety issues in deep learning solutions for this task. We 103 propose a knowledge unlearning-based approach, which refines pre-trained models by distinguishing 104 among three sets of training data and unlearning specific targets. This approach reduces ethical risks, 105 specifically the ability of deep learning models to enhance the properties of viruses, while maintaining the model's overall performance in designing normal, non-harmful proteins. Comprehensive 106 validation on multiple open benchmarks demonstrates the empirical significance of our proposed 107 PROEDIT compared to existing solutions and ablation models in terms of effectiveness, consistency,

and efficiency. Although this work is an initial exploration, we believe that safety concerns in AI for
 biology are critically important and merit greater attention and discussion.

111 112

113

114

122

123

140

141

2 PRELIMINARY: MODEL PRE-TRAINING AND KNOWLEDGE UNLEARNING

2.1 PRE-TRAINED PROTEIN LANGUAGE MODEL FOR MUTATION EFFECT PREDICTION

PLMs are the mainstream approach for learning protein sequence representations, including both BERT-style and GPT-style pre-training schemes. The former learns to recover masked tokens in the input sequence, and the latter generates tokens autoregressively. For mutation effect prediction, we implement the BERT-style approach. During training, a BERT-style model applies random masks to the input sequence. The training objective is to find θ , the optimal parameters that minimize the difference between the prediction of the masked amino acids (AAs) and the corresponding ground truth, *i.e.*,

$$\underset{\boldsymbol{\theta}}{\arg\min} \mathbb{E}_{\boldsymbol{x} \sim \boldsymbol{X}} \mathbb{E}_M - \sum_{i \in M} \log \mathrm{P}(\boldsymbol{x}_i | \boldsymbol{x}_{/M}; \boldsymbol{\theta}).$$
(1)

The conditional probability $P(x_i|x_{/M})$ of the *i*-th token x_i in the sequence is based on the unmasked part $x_{/M}$. The model learns to interpret the interactions of AAs within the protein sequence.

The trained model PLM°, obtained from (1), provides a summary matrix of the probability distribution for each AA in the sequence. This distribution has been shown to be effective for scoring mutation effects, especially when there is insufficient experimental data to support supervised learning. Given the AA probability distribution obtained from a pre-trained model PLM° for a wild-type protein, it can score relevant mutants of interest. Denote a |F|-site mutant by a set of triplets $F = \{(i, F_i, w_i) | i = 1, 2, ..., |F|\}$, where F_i and w_i are the residue types of the *i*th AA after and before the point mutation, respectively. The fitness score of the mutant F is:

score(
$$\boldsymbol{F}$$
) = $\sum_{i=1}^{|\boldsymbol{F}|} \log P(\boldsymbol{x}_i = \boldsymbol{F}_i | \boldsymbol{x}; PLM^{o}) - \log P(\boldsymbol{x}_i = \boldsymbol{w}_i | \boldsymbol{x}; PLM^{o}).$ (2)

The above zero-shot scoring function provides the *log-odds ratio* of mutants. Since most enzymes lack sufficient experimental labels to train a supervised learning model, this strategy is currently the most widely used scoring function in related research (Meier et al., 2021; Notin et al., 2024).

2.2 KNOWLEDGE UNLEARNING FOR PROTEIN LANGUAGE MODEL

Suppose an initial PLM^o is trained on a collection of protein sequences with arbitrary properties \mathbb{D}° . For a new input x, this well-trained model can provide the corresponding output $y = \text{PLM}^{\circ}(x)$, regardless of whether x and the associated property to modify is desired or undesired.

However, as mentioned in the previous section, protein engineering tasks desire a safe and responsible model to provide reliable enhancement strategies for normal proteins (such as industrial enzymes), while being incapable of modifying harmful proteins (such as viruses). Formally, if there is a set of desired normal proteins $(x^{norm}, y^{norm}) \in \mathbb{D}^{norm}$ and a set of undesired proteins $(x^{\text{fgt}}, y^{\text{fgt}}) \in \mathbb{D}^{\text{fgt}}$, we aim to learn a modified PLM^{new}, which reduces the effectiveness of PLM^o in understanding undesirable instances while maintaining its ability to infer desired instances, *i.e.*,

$$y^{\text{fgt}} \neq \text{PLM}^{\text{new}}(x^{\text{fgt}})$$

and $y^{\text{norm}} = \text{PLM}^{\text{new}}(x^{\text{norm}}).$ (3)

3 KNOWLEDGE UNLEARNING VIA MODEL RETRAINING

157 3.1 GENERAL OBJECTIVE

The overall goal of model unlearning in the context of mutation effect prediction is to reduce the model's ability to represent viruses while minimally affecting its representation capabilities for normal (non-virus) proteins. Methodologically, the aim is to update a pre-trained model PLM^o into a new model PLM^{new}, *i.e.*, updating the parameters from θ^{o} to θ . To achieve this, we prepare three

153 154 155

156

158

training datasets corresponding to three optimization objectives: the unlearning scope, the retention
 scope, and the corruption scope. The unlearned model PLM^{new} is expected to forget the knowledge
 in the unlearning scope while retaining the knowledge in the retention scope. Additionally, we de fine a corruption scope for virus-like proteins, on which we expect the unlearned model to maintain
 similar performance to the original model. Formally, we define the objective function as follows:

$$\underset{\boldsymbol{\theta}}{\operatorname{arg\,max}} \underbrace{-\mathbb{E}_{\boldsymbol{x} \sim \mathbb{D}^{\operatorname{fgt}}} \log \mathcal{P}_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{x}_{M})}_{\operatorname{Unlearn Scope}} + \underbrace{\mathbb{E}_{\boldsymbol{x} \sim \mathbb{D}^{\operatorname{norm}}} \log \mathcal{P}_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{x}_{M})}_{\operatorname{Retention Scope}} + \underbrace{\mathbb{E}_{\boldsymbol{x} \sim \mathbb{D}^{\operatorname{sim}}} \operatorname{KL} \left(\mathcal{P}_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{x}_{M}) \| \mathcal{P}_{\boldsymbol{\theta}^{\circ}}(\boldsymbol{x} | \boldsymbol{x}_{M}) \right)}_{\operatorname{Corruption Scope}}.$$
(4)

171 172 173

174

175 176

177

178

170

167 168 169

Unlearn Scope The first term, $-\mathbb{E}_{\boldsymbol{x}\sim\mathbb{D}^{\text{fgt}}} \log P_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{x}_M)$, measures the representation ability of PLM^{new} in recovering masked tokens in viruses. The model's parameters are updated based on a virus dataset \mathbb{D}^{fgt} to ensure the knowledge is forgotten. A model is considered to have effectively forgotten undesired virus knowledge if the recovery performance is poor, *i.e.*, if $-\mathbb{E}_{\boldsymbol{x}\sim\mathbb{D}^{\text{fgt}}} \log P_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{x}_M)$ is maximized.

179 **Retention Scope** The second term, $\mathbb{E}_{\boldsymbol{x} \sim \mathbb{D}^{norm}} \log P_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{x}_M)$, measures the effectiveness of PLM^{new} in recovering masked tokens in \boldsymbol{x}^{norm} , a set of non-virus normal proteins that are mutually exclusive to \mathbb{D}^{fgt} . A well-unlearned model is expected to minimize $\mathbb{E}_{\boldsymbol{x} \sim \mathbb{D}^{norm}} \log P_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{x}_M)$, indicating it retains knowledge relevant to normal proteins.

Corruption Scope The third term, $\mathbb{E}_{\boldsymbol{x} \sim \mathbb{D}^{\text{sim}}} \text{KL}(P_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{x}_M) || P_{\boldsymbol{\theta}^{\circ}}(\boldsymbol{x}|\boldsymbol{x}_M))$, focuses on a challenging subset of virus-like proteins $\mathbb{D}^{\text{sim}} \subset \mathbb{D}^{\text{norm}}$. As an augmentation, this term requires the unlearned model to minimize the difference (measured by the KL divergence) between PLM^{new} and the original model PLM^o, ensuring that forgetting viral knowledge does not disrupt the knowledge of virus-like normal proteins.

188 189

190

3.2 TRAINING SCHEME

191 **Data Preparation** To implement the untraining scheme, we divide the UniRef50 dataset ¹ into three sets. The first two sets, \mathbb{D}^{fgt} and \mathbb{D}^{norm} , are directly split from the processed UniRef50 dataset 192 based on their annotated Taxon IDs. Specifically, \mathbb{D}^{fgt} includes proteins whose Taxon IDs indicate 193 a biological lineage of viruses ². The remaining proteins form \mathbb{D}^{norm} . These two datasets, after pro-194 cessing, contain 65, 511, 306 and 564, 268 sequences, respectively. The statistics of these sequences 195 are detailed in Appendix A.1. For the virus-like proteins \mathbb{D}^{sim} , we extract them from \mathbb{D}^{norm} using 196 a retrieval module. Specifically, for each virus protein $x^{\text{fgt}} \in \mathbb{D}^{\text{fgt}}$, we pair the k-nearest proteins 197 from \mathbb{D}^{norm} based on the cosine similarity of their ESM-2 (650M) sequence embeddings. After preparing all three datasets, we conduct a random split on each of them, resulting in corresponding 199 training, validation, and test sets with a ratio of 8:1:1. 200

Model Optimization The trainable parameters θ of PLM^{new} are updated iteratively based on (4). To enhance training stability, we adopt an alternating micro-batch training strategy. Samples from the same batch originate exclusively from one of the datasets \mathbb{D}^{fgt} , \mathbb{D}^{norm} , or \mathbb{D}^{sim} , as defined in (5):

204 205

206

207 208

214

215

 $\mathcal{L}_{\text{batch}} = \begin{cases} -\frac{1}{|\mathcal{B}|} \sum_{\boldsymbol{x} \in \mathcal{B}} \log \mathcal{P}_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{x}_{M}), & \text{if } \mathcal{B} \subset \mathbb{D}^{\text{norm}} \\ \frac{1}{|\mathcal{B}|} \sum_{\boldsymbol{x} \in \mathcal{B}} \log \mathcal{P}_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{x}_{M}), & \text{if } \mathcal{B} \subset \mathbb{D}^{\text{fgt}} \\ \frac{1}{|\mathcal{B}|} \sum_{\boldsymbol{x} \in \mathcal{B}} \operatorname{KL}\left(\mathcal{P}_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{x}_{M}) \| \mathcal{P}_{\boldsymbol{\theta}^{\circ}}(\boldsymbol{x} | \boldsymbol{x}_{M})\right), & \text{if } \mathcal{B} \subset \mathbb{D}^{\text{sim}}. \end{cases}$ (5)

This approach ensures that the model focuses on a single objective at a time, thereby improving convergence and preventing interference between different learning objectives. The stopping criteria include four key considerations concerning perplexity and Spearman's ρ :

1. Perplexity of sampled data from \mathbb{D}^{norm} (smaller is better);

¹https://www.uniprot.org/help/downloads

²The biological lineage of a Taxon ID can be obtained from NCBI at https://www.ncbi.nlm.nih.gov/

- 216 2. Perplexity of sampled data from \mathbb{D}^{fgt} (larger is better);
 - 3. Spearman's ρ for assays of virus from ProteinGym (smaller is better);
 - 4. Spearman's ρ for assays of normal proteins from ProteinGym (larger is better).

We randomly sample 10% of instances as validation set. At the end of each epoch, we compute these metrics and the training will be stopped if any of the four metrics decreases for ten consecutive epochs.

225 3.3 Alternative Unlearning Methods

In addition to the optimization method we proposed above, strategies from other frameworks can also be adopted to unlearn PLMs. Below we brief four alternative strategies employed in unlearning LLMs. Their performance in unlearning PLMs will be compared in the subsequent section.

Gradient Ascent The first method uses gradient ascent (Tian et al., 2024) to forget learned knowledge. The learning objective remains consistent with the pre-training stage. When untraining a pre-trained MLM, this method trains on \mathbb{D}^{norm} and updates the model parameters by performing gradient ascent, *i.e.*, the opposite of descent, on (1).

Model Corruption with Random Labels The second approach is to fine-tune the model using randomly generated labels (Golatkar et al., 2020). Intuitively, by associating the data to be forgotten with random or incorrect labels, the model is expected to unlearn the associations it had previously made. In our case, this method trains using both \mathbb{D}^{fgt} and \mathbb{D}^{norm} . The ground truth labels from \mathbb{D}^{fgt} are randomly replaced with uniformly sampled tokens from the vocabulary, while the labels from \mathbb{D}^{norm} remain unchanged and are used to train with gradient descent.

240 241

218 219

220

221

222

223 224

229

234

Joint Gradient Ascent and Descent The third hybrid method leverages gradient ascent to forget undesired information and gradient descent to retain useful knowledge (Yao et al., 2023). By alternating between the two, the model aims to forget specific information while retaining as much overall performance as possible. We apply this strategy to untrain a PLM on both D^{fgt} and D^{norm}, using (1) as the training objective. Gradient ascent is applied to D^{fgt}, while gradient descent is applied to D^{norm}. It can be considered as a variant of PROEDIT that omits the corruption scope.

248 **Gradient Ascent with KL Constraint** The last strategy uses KL-divergence (Yao et al., 2023) to constrain the model's outputs, ensuring they do not stray too far from the original knowledge 249 during the unlearning process. This approach helps balance the unlearning task by maintaining a 250 good trade-off between forgetting and retaining information. The model is updated by performing 251 gradient ascent on \mathbb{D}^{fgt} and using KL divergence on \mathbb{D}^{norm} . Notably, the KL divergence on \mathbb{D}^{norm} 252 ensures that the model's outputs remain consistent with those of the original model. The key dif-253 ference between this method and the second method (joint gradient ascent and descent) is that the 254 former uses KL divergence to maintain output consistency, whereas the latter employs the MLM 255 pre-training objective to prevent forgetting knowledge of normal (non-virus) proteins.

- 256 257 258
- 4 EXPERIMENTS
- 259 260

4.1 EXPERIMENTAL PROTOCOL

261 **Setup** We compare the performance of PROEDIT and baseline methods on different benchmark 262 datasets. We use pre-trained ESM-2 (150M) and ESM-2 (650M) (Lin et al., 2023b) as the two base 263 models and apply PROEDIT for editing. To improve training efficiency, we limit each epoch to a 264 maximum of 2,000 batches, with each batch containing 4 samples. All trainable parameters are 265 updated using the ADAM (Kingma & Ba, 2015) optimizer with a learning rate of 1×10^{-5} . The 266 MLM pre-training objective remains consistent with that of ESM-2, as defined in (1). Specifically, 267 15% of the tokens in each sequence are selected for masking: 80% are replaced with the "[MASK]" token, 10% are substituted with a random token, and the remaining 10% are left unchanged. For the 268 four alternative unlearning methods introduced in Section 3.3, we adopt the same hyper-parameter 269 configurations as PROEDIT. All implementations are done using PyTorch (version 1.7.0), and the

Table 1: Statistic	s Summary on	the three	benchmarks
--------------------	--------------	-----------	------------

ľ	Name	Туре	# Assays (virus)	# Assays (non-virus)	# Mutations	# Train Samples	# Test Samples
Prot	teinGym	Zero-Shot	30	187	2,465,767	-	2,465,767
	AAV	Supervised	1	0	82,583	1,170	81,413
	GB1	Supervised	0	1	8,733	5,089	3,644

experiments are run on an NVIDIA[®] RTX 4090 GPU with 24GB VRAM, mounted on a server with Ubuntu 22.04 LTS operating system. All the details to reproduce our results are included in the submission, and the code will be made publicly available upon acceptance.

Benchmark Datasets We conduct a comprehensive evaluation of PROEDIT's editing capabilities, including both zero-shot and fine-tuning performance. The zero-shot prediction is evaluated on 217 deep mutational scanning (DMS) assays from **ProteinGym** (Notin et al., 2024), while fine-tuning is assessed on two supervised learning tasks: a viral dataset (**AAV**, Adeno-associated virus) and a non-viral dataset (**GB1**, binding domain of protein G, from Streptococcal bacteria) (Dallago et al., 2021). The statistics of these three benchmarks are summarized in Table 1. The training and evaluation procedures for these three tasks are as follows:

- 289 1. **ProteinGym** includes 30 viral assays and 187 non-viral assays. For each model, including 290 PROEDIT, we score mutants using (2) and calculate the Spearman's ρ correlation between the 291 predicted and experimental mutational scores. In evaluating model performance, we aim for a 292 high (closer to 1) Spearman's ρ on non-viral assays and a low Spearman's ρ on viral assays.
- 293 2. AAV is used for the first supervised learning task, which is a viral protein dataset originating from 294 the FLIP benchmark. We use the "one-vs-rest" split, which consists of 1, 170 single-order muta-295 tions for training and 81, 413 high-order mutations for testing. Model performance is evaluated 296 using Spearman's ρ , where a good model should have a lower Spearman's ρ .
- 3. GB1, a binding protein from Streptococcal bacteria, is used for the second supervised learning task. It is another dataset from the FLIP benchmark. We use the "low-vs-high" split from the FLIP benchmark, which includes 5, 089 low-fitness mutations for training and 3, 644 high-fitness mutations for testing. Similar to AAV, model performance is evaluated using Spearman's *ρ*. In this case, a good model should have a higher Spearman's *ρ*.
- For both supervised learning tasks, **GB1** and **AAV**, we added a regression MLP head to the model, which was fine-tuned on the respective training sets. The fine-tuning process used the following hyperparameters: the ADAMW (Kingma & Ba, 2015; Loshchilov et al., 2017) optimizer with a learning rate of 0.0005, a weight decay of 0.01, a batch size of 16, and a dropout rate of 0.1 for the output layer.
- 307 308 309

270

278

279

280 281

282

283

284

285

286

287

288

4.2 RESULTS ANALYSIS

We evaluate the performance of the proposed PROEDIT from three dimensions: effectiveness, con-310 sistency, and efficiency. The performance comparison of PROEDIT and baseline methods on mu-311 tation effect prediction tasks is reported in Table 2. We use ESM-2 (650M) and ESM-2 (150M) as 312 our two base models. We also compare several other pre-trained PLM models, including MIF-ST 313 Yang et al. (2023), CARP Yang et al. (2024), and ESM-1v Meier et al. (2021). Additionally, we 314 include a vanilla Transformer Vaswani et al. (2017) for the two supervised learning tasks. The archi-315 tecture of vanilla Transformer is the same as ESM-2 and the only difference is that the parameters 316 of the vanilla Transformer are randomly initialized before training. The initialization configuration 317 is according to ESM-2 and the random seed is 42.

318

Effectiveness. From the results reported in Table 2, it can be observed that PROEDIT significantly
 reduces the performance of the base model on ProteinGym-virus and AAV, while maintaining
 performance on ProteinGym-non-virus and GB1. This indicates that PROEDIT effectively assists
 the pre-trained base model in unlearning viral knowledge in both zero-shot and fine-tuning tasks.
 Specifically, for PROEDIT (650M), compared to ESM-2 (650M), the Spearman's ρ correlation on
 ProteinGym-virus dropped from 0.24 to 0.08, retaining only 30% of the original performance. On

		zero-shot prediction		fine-tuning	
Model	version	$\hline \textbf{ProteinGym} (virus) \downarrow \\$	ProteinGym (non-virus) ↑	$\overline{\mathbf{AAV}\downarrow}$	GB1↑
Transformer	vanilla	-	-	0.30	0.08
MIF-ST	-	0.40	0.43	-	0.22
CARP	640M	0.27	0.37	0.43	0.48
ESM-1v	-	0.28	0.44	0.37	0.27
ESM-2	150M	0.13	0.45	0.08	0.15
ESM-2	650M	0.24	0.48	0.35	0.17
ProEdit	150M	0.08	0.43	-0.16	0.13
ProEdit	650M	0.07	0.47	-0.18	0.24

Table 2: Spearman's ρ correlation of mutation effect prediction by different methods.

[†] The top two are highlighted by First and Second.



Figure 3: Individual prediction of assays in ProteinGym (virus) by PROEDIT and ESM-2.



Figure 4: (a) Validation curves and (b) UMAP embedding of PROEDIT (650M) and ESM-2 (650M).

the viral protein AAV, PROEDIT (650M) dropped to -0.18 and PROEDIT (150M) dropped to -0.16, which are significantly lower than the original scores at 0.34 and 0.08, respectively. The detailed prediction performance on individual assays in **ProteinGym** (virus) is visualized in Figure 3, where PROEDIT greatly reduces the Spearman's ρ on the majority of the assays in comparison to the performance of the base model (ESM-2-650M), with specific numbers provided in Appendix A.2. Addi-tionally, comparing the results in Figures 4(b)-(c) shows that PROEDIT successfully scrambled viral information after unlearning. Here, we randomly selected 1,000 virus and 1,000 non-viral protein sequences and extracted hidden representations using both PROEDIT (650M) and ESM-2 (650M), followed by dimension reduction with UMAP (McInnes et al., 2018). It is evident that, compared to ESM-2, the representations of virus and non-viral proteins encoded by PROEDIT are more indistinguishable, further indicating the successful unlearning of viral proteins by PROEDIT. Similarly, Figure 4(a) displays the validation curves for PROEDIT (650M) and ESM-2 (650M) trained on AAV (virus) in a supervised learning setup. It shows clearly that the performance of PROEDIT remains constantly at a low level during training on this virus dataset, demonstrating the effectiveness of the unlearning process.

Consistency When the model forgets viral knowledge, it should minimize any negative impact on general protein knowledge. Specifically, the performance reduction in scoring non-viral proteins be-



378 Table 3: Zero-shot prediction performance on ProteinGym by PROEDIT and alternative unlearning 379 methods.

380

scores in Table 2 for ProteinGym (non-virus) and GB1. On the ProteinGym (non-virus), the score of PROEDIT (650M) decreased by only 0.01 compared to ESM-2 (650M), and PROEDIT (150M) decreased by merely 0.02 compared to ESM-2 (150M). Furthermore, on GB1, comparing the same parameter versions of ESM-2 and PROEDIT reveals that the unlearned model achieves better performance on this non-viral protein than the base model, surpassing all other baseline models. This is evident that the unlearning method is capable of retaining general protein knowledge when removing

Efficiency We show that editing pre-trained PLMs with PROEDIT would not introduce significant additional computational costs. In our experiments, training PROEDIT costs 7.27 GPU hours for updating parameters in ESM-2 (650M), and 5.25 GPU hours for updating parameters in ESM-2 (150M). 426

428 4.3 ADDITIONAL INVESTIGATIONS

429

427

Alternative Unlearning Methods Table 3 compares the zero-shot prediction performance of four 430 alternative unlearning strategies mentioned in Section 3.3 on ProteinGym. Although all five un-431 learning methods (including ours) can assist a pre-trained PLM in forgetting harmful knowledge

about viruses, the alternative methods significantly reduce performance on non-viral proteins at the
same time. In contrast, our PROEDIT is more practically meaningful. It is capable of maintaining
both effectiveness and consistency, thereby enhancing model safety while preserving reliability on
general protein tasks.

436

437 **Construction of the Virus-Like Dataset** Figure 5 explores the impact of different choices on k in the retrieval module when constructing \mathbb{D}^{sim} on the results. Panels (a) and (b) report the changes 438 in Spearman's ρ and perplexity as k varies, with k = 0 indicating the scores without editing. As k 439 increases gradually from 1 to 5, the model's performance on non-viral proteins slowly declines. A 440 possible explanation is that increasing k would directly enlarge \mathbb{D}^{sim} , causing some normal proteins 441 that are not very similar to viruses to be included in the corruption scope, thus interfering with 442 the model parameter updates. Additionally, in **ProteinGym** (virus), increasing k does not help the 443 model forget viral information more effectively. Therefore, in practice, we set k = 1 to perform an 444 effective and efficient unlearning scheme. 445

446 **Change of Parameters Before and After Unlearning** Figure 6 further investigates the changes in 447 learnable parameters of the Transformer layers before and after editing to provide additional insights 448 on the overall impact of the unlearning module. Taking PROEDIT (650M) and ESM-2 (650M) as 449 examples, Figure 6(a) displays the ten network layers with the largest and smallest differences in L2 450 norm. The specific L2 norm of all parameters of the Transformer layers are provided in Appendix 451 A.3. Overall, among the 33 Transformer layers, parameters of layers closer to the output show greater changes, while layers closer to the input exhibit smaller changes. This indicates that the 452 model's forgetting primarily occurs in the last few layers of the Transformer. Figure 6(b) illustrates 453 the parameter changes of the layers with the largest and smallest variations. Despite the differing 454 magnitudes of change, both exhibit a symmetric bell-shaped distribution. 455

456

5 RELATED WORK

457 5 F 458

Pre-trained Protein Language Models Analogous to NLP models, PLMs typically treat AAs 459 as tokens and use Transformer-based layers (Vaswani et al., 2017) to analyze co-evolutionary in-460 formation from millions to billions of protein sequences and summarize vector representations for 461 sequences. Pre-trained PLMs can be categorized into three types. The most common are encoder-462 only models, which follow the BERT framework (Devlin et al., 2018) and train the model to recover 463 randomly masked AA types (Meier et al., 2021; Rives et al., 2021; Elnaggar et al., 2021; Tan et al., 464 2024b). Decoder-only models, in comparison, are trained to optimize next-token prediction, which 465 is frequently used for sequence design (Ferruz et al., 2022; Notin et al., 2022a; Madani et al., 2023). 466 Other models adopt hybrid encoder-decoder architectures to learn outputs that are sufficiently simi-467 lar to the input sequences (Du et al., 2022; Elnaggar et al., 2023; Heinzinger et al., 2023).

468 469

481

Mutation Effect Prediction When applying models to enhance protein properties and function-470 alities, some studies adopt a "pre-training then fine-tuning" approach to enhance the model's un-471 derstanding of a particular protein assay, such as using existing experimental data for supervised learning (Li et al., 2023; Zhou et al., 2024c; Tan et al., 2024a) or incorporating homologous se-472 quences during training (Rao et al., 2021; Notin et al., 2022b). However, due to the lack of publicly 473 available mutation effect data for most proteins and the variability among assays, the mainstream 474 approach remains zero-shot methods for mutation effect prediction. Considering the pivotal role of 475 structure in determining a protein's function, many recent methods integrate geometric deep learning 476 methods (Lu et al., 2022; Tan et al., 2023; Zhou et al., 2024a; Tan et al., 2024c) or extract structure 477 tokens (Su et al., 2023; Li et al., 2024a) to enhance the local interaction of spatially connected AAs. 478 With the introduction of large-scale deep mutational scanning benchmark datasets like ProteinGym 479 (Notin et al., 2024), an increasing number of models have been developed and extensively validated 480 on a wide range of protein assays to demonstrate their effectiveness and generalizability.

Knowledge Unlearning and AI Safety As emerging models grow larger and training data be comes more diverse, increasing attention is being directed toward developing approximate unlearn ing algorithms, such as data-reversed training (Chundawat et al., 2023) and optimization-based un learning (Guo et al., 2020; Neel et al., 2021). In the NLP field, particularly with LLMs, AI safety has
 caught increasing attention. Knowledge unlearning, in this context, trains models to reject sensitive

responses (Yu et al., 2023; Yao et al., 2023; Tian et al., 2024; Li et al., 2024b). Similar discussions have emerged in other fields such as Computer Vision (Kim & Woo, 2022; Lin et al., 2023a; Tarun et al., 2023). However, in AI for biology, particularly in protein engineering, the safety and responsibility of developed deep learning models remain under-explored.

490 491 492

6 CONCLUSION AND DISCUSSION

493 This study addresses a critical safety concern in mutation effect prediction, a core task in protein en-494 gineering, by proposing a novel knowledge unlearning-based framework, PROEDIT. Our approach 495 enables pre-trained PLMs to selectively forget virus-related information while preserving their ca-496 pacity to predict and design non-viral proteins. Through empirical validation on multiple bench-497 marks, we demonstrated that PROEDIT effectively reduces the risk of enhancing viral properties 498 without compromising the performance of beneficial proteins. This contributes to the growing need 499 for ethical and responsible AI in scientific applications, particularly in biosafety-sensitive domains 500 like protein engineering.

501 The rapid development of deep learning techniques in recent years has led to an increasing number 502 of powerful solutions to biological challenges. The growing attention to these advancements has 503 significantly driven improvements in the prediction and generation performance of biological en-504 tities, such as drug discovery, enzyme engineering, and protein design. However, alongside these 505 technological advancements, we emphasize that ensuring their ethical and responsible use is equally crucial. We hope this work inspires more researchers to explore safety concerns in AI-driven pro-506 tein engineering and to extend the unlearning framework to other safety-critical applications. By 507 developing models that excel in predictive power while also addressing potential risks, the scientific 508 community can promote safer and more responsible advancements in AI-driven biological research. 509

510 511

520

521

522

References

- Vikram S Chundawat, Ayush K Tarun, Murari Mandal, and Mohan Kankanhalli. Zero-shot machine unlearning. *IEEE Transactions on Information Forensics and Security*, 18:2345–2354, 2023.
- Christian Dallago, Jody Mou, Kadina E Johnston, Bruce Wittmann, Nick Bhattacharya, Samuel Goldman, Ali Madani, and Kevin K Yang. FLIP: Benchmark tasks in fitness landscape inference for proteins. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. URL https://openreview.net/forum?id= p2dMLEwL8tF.
 - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv:1810.04805*, 2018.
- 523 Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. GLM:
 524 General language model pretraining with autoregressive blank infilling. In *Proceedings of the*525 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers),
 526 pp. 320–335, 2022.
- Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Yu Wang, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, et al. Prottrans: Toward understanding the language of life through self-supervised learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):7112–7127, 2021.
- Ahmed Elnaggar, Hazem Essam, Wafaa Salah-Eldin, Walid Moustafa, Mohamed Elkerdawy, Charlotte Rochereau, and Burkhard Rost. Ankh: Optimized protein language model unlocks general-purpose modelling. *arXiv:2301.06568*, 2023.
- Noelia Ferruz, Steffen Schmidt, and Birte Höcker. ProtGPT2 is a deep unsupervised language model
 for protein design. *Nature Communications*, 13(1):4348, 2022.
- Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net:
 Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9304–9312, 2020.

555

582

583

540	Aditya Golatkar, Alessandro Achille, Avinash Ravichandran, Marzia Polito, and Stefano Soatto.
541	Mixed-privacy forgetting in deep networks. In Proceedings of the IEEE/CVF conference on com-
542	nuter vision and pattern recognition. pp. 792–801. 2021.
543	

- Chuan Guo, Tom Goldstein, Awni Hannun, and Laurens Van Der Maaten. Certified data removal
 from machine learning models. In *International Conference on Machine Learning*, pp. 3832–3842. PMLR, 2020.
- Haohuai He, Bing He, Lei Guan, Yu Zhao, Feng Jiang, Guanxing Chen, Qingge Zhu, Calvin YuChian Chen, Ting Li, and Jianhua Yao. De novo generation of sars-cov-2 antibody cdrh3 with a
 pre-trained generative large language model. *Nature Communications*, 15(1):6867, 2024.
- Michael Heinzinger, Konstantin Weissenow, Joaquin Gomez Sanchez, Adrian Henkel, Martin Steinegger, and Burkhard Rost. ProstT5: Bilingual language model for protein sequence and structure. *bioRxiv*, pp. 2023–07, 2023.
 - Junyaup Kim and Simon S Woo. Efficient two-stage model retraining for machine unlearning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4361–4369, 2022.
- Diederik P Kingma and Jimmy Ba. ADAM: A method for stochastic optimization. In *Proceedings* of International Conference on Learning Representation, 2015.
- 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, 2019.
- Mingchen Li, Liqi Kang, Yi Xiong, Yu Guang Wang, Guisheng Fan, Pan Tan, and Liang Hong.
 Sesnet: sequence-structure feature-integrated deep learning method for data-efficient protein engineering. *Journal of Cheminformatics*, 15(1):12, 2023.
- Mingchen Li, Yang Tan, Xinzhu Ma, Bozitao Zhong, Huiqun Yu, Ziyi Zhou, Wanli Ouyang, Bingxin Zhou, Liang Hong, and Pan Tan. Prosst: Protein language modeling with quantized structure and disentangled attention. *bioRxiv*, pp. 2024–04, 2024a.
- 568 Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D
 569 Li, Ann-Kathrin Dombrowski, Shashwat Goel, Gabriel Mukobi, et al. The wmdp benchmark:
 570 Measuring and reducing malicious use with unlearning. In *Forty-first International Conference*571 *on Machine Learning*, 2024b.
- Shen Lin, Xiaoyu Zhang, Chenyang Chen, Xiaofeng Chen, and Willy Susilo. Erm-ktp: Knowledge-level machine unlearning via knowledge transfer. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20147–20155, 2023a.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,
 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level
 protein structure with a language model. *Science*, 379(6637):1123–1130, 2023b.
- Shiwei Liu, Tian Zhu, Milong Ren, Chungong Yu, Dongbo Bu, and Haicang Zhang. Predicting
 mutational effects on protein-protein binding via a side-chain diffusion probabilistic model. *Ad- vances in Neural Information Processing Systems*, 36, 2024.
 - Ilya Loshchilov, Frank Hutter, et al. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101*, 5, 2017.
- Hongyuan Lu, Daniel J Diaz, Natalie J Czarnecki, Congzhi Zhu, Wantae Kim, Raghav Shroff,
 Daniel J Acosta, Bradley R Alexander, Hannah O Cole, Yan Zhang, et al. Machine learning-aided
 engineering of hydrolases for PET depolymerization. *Nature*, 604(7907):662–667, 2022.
- Ali Madani, Ben Krause, Eric R Greene, Subu Subramanian, Benjamin P Mohr, James M Holton, Jose Luis Olmos, Caiming Xiong, Zachary Z Sun, Richard Socher, et al. Large language models generate functional protein sequences across diverse families. *Nature Biotechnology*, 41(8):1099–1106, 2023.
- Leland McInnes, John Healy, Nathaniel Saul, and Lukas Grossberger. Umap: Uniform manifold approximation and projection. *The Journal of Open Source Software*, 3(29):861, 2018.

619

626

633

594	Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, and Alex Rives, Language
595	models enable zero-shot prediction of the effects of mutations on protein function. Advances in
596	neural information processing systems, 34:29287–29303, 2021.
597	J. I. I. G. J. I. I. J. I. J.

- Harini Narayanan, Fabian Dingfelder, Alessandro Butté, Nikolai Lorenzen, Michael Sokolov, and
 Paolo Arosio. Machine learning for biologics: opportunities for protein engineering, developabil ity, and formulation. *Trends in Pharmacological Sciences*, 42(3):151–165, 2021.
- Seth Neel, Aaron Roth, and Saeed Sharifi-Malvajerdi. Descent-to-delete: Gradient-based methods for machine unlearning. In *Algorithmic Learning Theory*, pp. 931–962. PMLR, 2021.
- Pascal Notin, Mafalda Dias, Jonathan Frazer, Javier Marchena Hurtado, Aidan N Gomez, Debora
 Marks, and Yarin Gal. Tranception: protein fitness prediction with autoregressive transformers and inference-time retrieval. In *International Conference on Machine Learning*, pp. 16990–
 17017. PMLR, 2022a.
- Pascal Notin, Lood Van Niekerk, Aaron W Kollasch, Daniel Ritter, Yarin Gal, and Debora Susan
 Marks. Trancepteve: Combining family-specific and family-agnostic models of protein sequences
 for improved fitness prediction. 2022b.
- Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood Van Niekerk, Steffanie Paul, Han Spinner, Nathan Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, et al. Proteingym: large-scale benchmarks for protein fitness prediction and design. In *Advances in Neural Information Processing Systems*, volume 36, 2024.
- Jeffrey Ouyang-Zhang, Daniel Diaz, Adam Klivans, and Philipp Krähenbühl. Predicting a protein's stability under a million mutations. *Advances in Neural Information Processing Systems*, 36, 2024.
- Fabrizio Pucci, Martin Schwersensky, and Marianne Rooman. Artificial intelligence challenges for
 predicting the impact of mutations on protein stability. *Current Opinion in Structural Biology*, 72:161–168, 2022.
- Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu, and Alexander Rives. MSA transformer. In *International Conference on Machine Learning*, pp. 8844–8856. PMLR, 2021.
- Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo,
 Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from
 scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15):e2016239118, 2021.
- Anton Sinitsin, Vsevolod Plokhotnyuk, Dmitry Pyrkin, Sergei Popov, and Artem Babenko. Editable
 neural networks. In *International Conference on Learning Representations*, 2019.
- Jin Su, Chenchen Han, Yuyang Zhou, Junjie Shan, Xibin Zhou, and Fajie Yuan. Saprot: protein language modeling with structure-aware vocabulary. In *The Twelfth International Conference on Learning Representations*, 2023.
- Yang Tan, Bingxin Zhou, Lirong Zheng, Guisheng Fan, and Liang Hong. Semantical and topological
 protein encoding toward enhanced bioactivity and thermostability. *bioRxiv*, pp. 2023–12, 2023.
- Yang Tan, Mingchen Li, Bingxin Zhou, Bozitao Zhong, Lirong Zheng, Pan Tan, Ziyi Zhou, Huiqun
 Yu, Guisheng Fan, and Liang Hong. Simple, efficient, and scalable structure-aware adapter boosts
 protein language models. *Journal of Chemical Information and Modeling*, 2024a.
- Yang Tan, Mingchen Li, Ziyi Zhou, Pan Tan, Huiqun Yu, Guisheng Fan, and Liang Hong. Peta:
 evaluating the impact of protein transfer learning with sub-word tokenization on downstream applications. *Journal of Cheminformatics*, 16(1):92, 2024b.
- 647 Yang Tan, Jia Zheng, Liang Hong, and Bingxin Zhou. Protsolm: Protein solubility prediction with multi-modal features. *arXiv:2406.19744*, 2024c.

648 Ayush K Tarun, Vikram S Chundawat, Murari Mandal, and Mohan Kankanhalli. Fast yet effective 649 machine unlearning. IEEE Transactions on Neural Networks and Learning Systems, 2023. 650 Bozhong Tian, Xiaozhuan Liang, Siyuan Cheng, Oingbin Liu, Mengru Wang, Dianbo Sui, Xi Chen, 651 Huajun Chen, and Ningyu Zhang. To forget or not? towards practical knowledge unlearning for 652 large language models. arXiv:2407.01920, 2024. 653 654 Timothy Truong Jr and Tristan Bepler. Poet: A generative model of protein families as sequences-655 of-sequences. Advances in Neural Information Processing Systems, 36:77379–77415, 2023. 656 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, 657 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in Neural Informa-658 tion Processing Systems, 30, 2017. 659 Meng Wang, Jonathan Patsenker, Henry Li, Yuval Kluger, and Steven H Kleinstein. Supervised fine-661 tuning of pre-trained antibody language models improves antigen specificity prediction. bioRxiv, 662 2024a. 663 Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, et al. Knowledge editing for 664 large language models: A survey. arXiv:2310.16218, 2023. 665 666 Weiqi Wang, Zhiyi Tian, and Shui Yu. Machine unlearning: A comprehensive survey. 667 arXiv:2405.07406, 2024b. 668 Kevin K Yang, Niccolò Zanichelli, and Hugh Yeh. Masked inverse folding with sequence transfer for 669 protein representation learning. Protein Engineering, Design and Selection, 36:gzad015, 2023. 670 671 Kevin K Yang, Nicolo Fusi, and Alex X Lu. Convolutions are competitive with transformers for 672 protein sequence pretraining. *Cell Systems*, 15(3):286–294, 2024. 673 Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. In NeurIPS2023 674 Workshop Socially Responsible Language Modelling Research, 2023. URL https:// 675 openreview.net/forum?id=wKe6jE065x. 676 677 Charles Yu, Sullam Jeoung, Anish Kasi, Pengfei Yu, and Heng Ji. Unlearning bias in language models by partitioning gradients. In Findings of the Association for Computational Linguistics: 678 ACL 2023, pp. 6032–6048, 2023. 679 680 Bingxin Zhou, Lirong Zheng, Banghao Wu, Yang Tan, Outongyi Lv, Kai Yi, Guisheng Fan, and 681 Liang Hong. Protein engineering with lightweight graph denoising neural networks. Journal of 682 Chemical Information and Modeling, 2024a. 683 Bingxin Zhou, Lirong Zheng, Banghao Wu, Kai Yi, Bozitao Zhong, Yang Tan, Qian Liu, Pietro 684 Liò, and Liang Hong. A conditional protein diffusion model generates artificial programmable 685 endonuclease sequences with enhanced activity. *Cell Discovery*, 10(1):95, 2024b. 686 687 Ziyi Zhou, Liang Zhang, Yuanxi Yu, Banghao Wu, Mingchen Li, Liang Hong, and Pan Tan. Enhanc-688 ing efficiency of protein language models with minimal wet-lab data through few-shot learning. 689 Nature Communications, 15(1):5566, 2024c. 690 691 692 693 694 696 697 699 700

A APPENDIX

702

703 704 705

706 707

708 709 710

711

712

A.1 DATASET STATISTICS

0.03

0.02

0.01

0.00

М

С В Н Т

V Y

745 746

747

748 749

750 751

752 753 This section presents statistics of the viral sequences and non-viral sequences in **UniRef50**, including the number of sequences, amino acid distribution, sequence lengths, and more.

Virus Sequences There are a total of **564,268** sequences. The length distribution and amino acid distribution of them are shown in Table 7 and Table 8, respectively.





QRA

14

I N E Amino Acids

Figure 8: Amino acid distribution of virus sequences in UniRef50.

LKGSF

PWXBZ







Table 4: Spearman's ρ correlation of mutation effect prediction by PROEDIT in ProteinGym (virus).

878			
879	DMS_ID	ProEdit (650M)	ESM-2 (650M)
880	A0A140D2T1_ZIKV_Sourisseau_2019	0.084	0.071
881	A0A192B1T2_9HIV1_Haddox_2018	-0.037	0.064
882	A0A2Z5U3Z0_9INFA_Doud_2016	-0.047	0.492
883	A0A2Z5U3Z0_9INFA_Wu_2014	0.039	0.452
оол	A4D664_9INFA_Soh_2019	0.118	0.137
004	C6KNH7_9INFA_Lee_2018	0.041	0.464
885	CAPSD_AAV2S_Sinai_2021	-0.127	0.200
886	ENV_HV1B9_DuenasDecamp_2016	0.006	0.042
887	ENV_HV1BR_Haddox_2016	-0.015	0.036
888	HCP_LAMBD_Tsuboyama_2023_2L6Q	0.471	0.695
889	I6TAH8_I68A0_Doud_2015	0.023	0.017
890	NCAP_I34A1_Doud_2015	0.070	0.020
891	NRAM_I33A0_Jiang_2016	0.040	0.166
892	PA_I34A1_Wu_2015	0.079	0.038
893	POLG_CXB3N_Mattenberger_2021	0.092	0.349
894	POLG_DEN26_Suphatrakul_2023	0.105	0.143
805	POLG_HCVJF_Qi_2014	0.106	0.127
095	POLG_PESV_Tsuboyama_2023_2MXD	0.100	0.406
090	Q2N0S5_9HIV1_Haddox_2018	-0.037	0.028
897	R1AB_SARS2_Flynn_2022	0.003	0.118
898	RDRP_I33A0_Li_2023	0.038	0.290
899	REV_HV1H2_Fernandes_2016	0.133	0.227
900	RPC1_BP434_Tsuboyama_2023_1R69	0.534	0.705
901	SPIKE_SARS2_Starr_2020_binding	-0.129	-0.015
902	SPIKE_SARS2_Starr_2020_expression	-0.057	0.030
903	TAT_HV1BR_Fernandes_2016	-0.081	-0.045
904	VG08_BPP22_Tsuboyama_2023_2GP8	0.450	0.662
905	VRPI_BPT7_Tsuboyama_2023_2WNM	0.029	0.576

Table 5: Parameter differences' L2 norm between the hidden and output of transformer layers of ESM-2 (650M) and PROEDIT (650M).

		1	
Parameter	L2 Norm	Parameter	L2 Norm
32 (output)	1.916	20 (hidden)	0.934
prediction head	1.724	21 (output)	0.933
32 (hidden)	1.458	19 (hidden)	0.932
31 (output)	1.348	7 (hidden)	0.932
31 (hidden)	1.260	11 (output)	0.932
30 (output)	1.140	10 (output)	0.931
30 (hidden)	1.105	9 (output)	0.928
29 (output)	1.058	17 (hidden)	0.928
28 (output)	1.048	20 (output)	0.928
29 (hidden)	1.030	19 (output)	0.926
28 (hidden)	1.005	18 (hidden)	0.925
10 (hidden)	0.993	6 (hidden)	0.923
9 (hidden)	0.984	13 (output)	0.918
11 (hidden)	0.982	12 (output)	0.918
27 (hidden)	0.981	17 (output)	0.914
27 (output)	0.976	18 (output)	0.914
26 (hidden)	0.965	16 (output)	0.912
26 (output)	0.960	15 (output)	0.911
25 (hidden)	0.956	14 (output)	0.910
12 (hidden)	0.955	7 (output)	0.907
24 (hidden)	0.948	8 (output)	0.904
23 (hidden)	0.947	5 (hidden)	0.904
25 (output)	0.945	6 (output)	0.898
13 (hidden)	0.945	5 (output)	0.888
8 (hidden)	0.945	4 (hidden)	0.863
22 (hidden)	0.943	4 (output)	0.855
24 (output)	0.943	3 (output)	0.748
23 (output)	0.941	3 (hidden)	0.745
22 (output)	0.940	2 (output)	0.720
15 (hidden)	0.938	2 (hidden)	0.719
14 (hidden)	0.938	1 (hidden)	0.603
21 (hidden)	0.937	1 (output)	0.596
16 (hidden)	0.935	0 (output)	0.560
		0 (hidden)	0.544