

# CoPeP: BENCHMARKING CONTINUAL PRETRAINING FOR PROTEIN LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In recent years, protein language models (pLMs) have gained significant attention for their ability to capture the structure and function of proteins, accelerating the discovery of new therapeutic drugs. These models are typically trained on large, evolving corpora of proteins that are continuously updated by the biology community. The dynamic nature of these datasets motivates the need for continual learning, not only to keep up with the ever-growing dataset sizes, but also as an opportunity to take advantage of the temporal meta-information that is created during this process. As a result, we introduce the Continual Pretraining of Protein Language Models (CoPeP) benchmark, a novel benchmark for evaluating continual learning approaches on pLMs. Specifically, we curate a sequence of protein datasets from the UniProt database spanning 10 years and define metrics to assess the performance of pLMs on diverse protein understanding tasks. We evaluate several methods from the continual learning literature, including replay, unlearning, and plasticity-based methods, some of which have never been applied to models and data of this scale. Our findings reveal that incorporating temporal meta-information improves perplexity by up to 20% over training from scratch on the latest snapshot of the database, and that several continual learning-based methods outperform naive continual pretraining. The CoPeP benchmark presents an exciting opportunity for studying these methods at scale on an impactful, real-world application.

## 1 INTRODUCTION

Proteins are the fundamental building blocks of life, acting as the primary machinery of all living organisms. Their function is mostly determined by their three-dimensional shape, which in turn is encoded into a linear sequence of 20 distinct amino acids. Predicting the properties of a protein from its sequence is one of the core challenges in computational biology. Recently, protein language models (pLMs) have emerged as an effective and scalable solution (Rives et al., 2021; Lin et al., 2023; Madani et al., 2023; Nijkamp et al., 2023; Fournier et al., 2024). By treating proteins as a language where amino acids are the “letters”, assembling into regions as “words”, themselves assembling into whole proteins as “sentences”, pLMs can discover the relationship between sequence, structure, and function from large databases (Rives et al., 2021; Notin et al., 2023). This allows them to accurately predict a protein’s properties and even to design new proteins for specific applications (Hayes et al., 2025), greatly accelerating drug discovery.

Despite their effectiveness, pLMs face a significant challenge in the dynamic nature of their training data (Fournier et al., 2024). These models are typically trained on enormous, ever-expanding public databases like the UniProt Knowledgebase (The UniProt Consortium, 2025), which are continuously updated. Each year, millions of new protein sequences are deposited, and millions of others are curated out after being identified as non-proteins. Consequently, the practice of retraining models from scratch on each new data release is becoming computationally prohibitive. This challenge, however, also presents a unique opportunity. The temporal evolution of these databases provides valuable metadata. Sequences that persist over time serve as strong examples of true proteins, while those that are later curated out can be treated as implicit examples of likely non-proteins. By leveraging this history, a model can more effectively learn the language of proteins.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

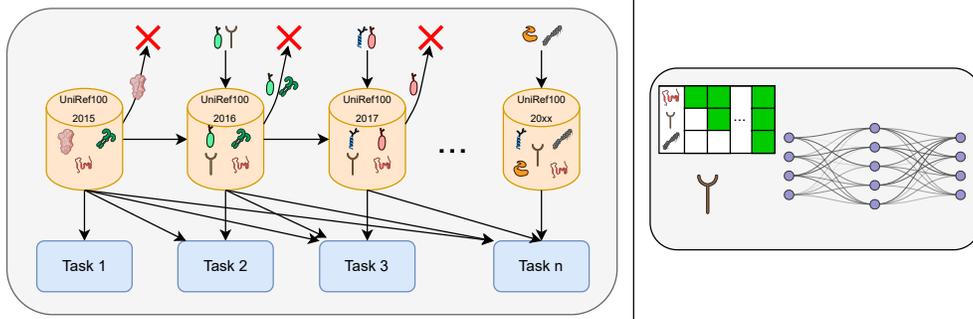


Figure 1: **Left:** The benchmark data process. For every year, we pull the latest UniRef100 release, which reflects the continuous discovery and curation of proteins by biologists. Each task in our benchmark is cumulative, containing all proteins from the start date up to a given year. **Right:** The model training setup. A pLM is trained on all protein sequences and can optionally leverage temporal metadata, such as the number of releases each protein has been a part of, which can be used by methods such as temporally weighted replay, model unlearning.

Although continual learning is a well-established field (Wang et al., 2024) with many artificial benchmarks, there is a growing demand for more realistic alternatives. Indeed, while these controlled environments are perfect for measuring loss of plasticity and catastrophic forgetting, they do not reflect the scale and complexity of real-world data. With the rise of Large Language Models (LLMs), there has been much interest in the community to explore ways to continually update these models with new information. One approach involves limited unlearning or updating a small number of facts that the model might have memorized (Bourtole et al., 2021; Yao et al., 2023). Other works try to extend the pretraining process itself with new datasets (Gupta et al., 2023; Abbes et al., 2025; Ke et al., 2022; Li et al., 2025). Despite this interest, however, there are not many general-purpose continual pretraining datasets where the goal is to extend the pretraining phase, and most academic works end up using domain-adaptive pretraining setups (Ke et al., 2022; Yıldız et al., 2025).

To bridge this gap, we introduce the Continual Pretraining for Protein Language Models (CoPeP) benchmark. CoPeP provides a realistic, large-scale solution for evaluating continual learning approaches on pLMs. We curate a sequence of protein datasets from 8 yearly releases of the UniProt database, giving us a unique opportunity to study how models adapt to continuously evolving data. We evaluate several state-of-the-art methods from the continual learning literature, including Gradient Ascent (Golotkar et al., 2020), Hare Tortoise (Lee et al., 2024), Replay (Rolnick et al., 2019; Chaudhry et al., 2019), and Shrink and Perturb (Ash & Adams, 2020), applying some of them for the first time at a scale comparable to real-world applications. We evaluate our models on two types of tasks: (1) a high-quality validation set of experimentally verified proteins to assess performance on natural protein distributions (Fournier et al., 2024); (2) ProteinGym (Notin et al., 2023), which benchmarks the ability to predict the effects of protein mutations. Our findings reveal that several of these methods improve performance over naive continual pretraining, and that leveraging temporal metadata yields a measurable improvement over models trained on individual years.

Our contributions are three-fold. First, we introduce CoPeP, a new benchmark for continual learning on real-world protein databases. Second, we are the first to apply and evaluate several state-of-the-art continual learning methods on a problem of this scale and complexity. Finally, we demonstrate that temporal metadata contained in the history of proteins being added or removed from the database can be leveraged to improve the performance of pLMs beyond that of single years.

## 2 RELATED WORK

**Continual Learning and Model Updating** Continual learning is a machine learning paradigm in which models are trained incrementally on a sequence of data or tasks, aiming to accumulate and update knowledge continuously much like humans do. Research in this area primarily focuses on two key challenges: catastrophic forgetting, the loss of previously acquired knowledge (McCloskey

108 & Cohen, 1989; Kirkpatrick et al., 2017), and loss of plasticity, the reduced ability to adapt to new  
109 data (Dohare et al., 2024). While some studies investigate continual learning under natural data  
110 shifts (Koh et al., 2021; Lin et al., 2021; Cai et al., 2021; Bornschein et al., 2023), the datasets  
111 used are typically much smaller than modern pretraining corpora. Most of the research in continual  
112 learning considers smaller academic datasets like CIFAR-10 and MNIST (Goodfellow et al., 2013;  
113 Zenke et al., 2017; Krizhevsky et al., 2009; Rebuffi et al., 2017) that allow for controlled experimen-  
114 tal setups and the study of severe distribution shifts that may be rare in natural data. However, the  
115 limited scale of these datasets raises questions about how well existing methods generalize to larger  
116 and more complex scenarios.

117 More recently, the field has started to shift toward updating large pretrained models. This includes  
118 model editing, which updates specific facts in the model without full retraining (Meng et al., 2022;  
119 Mitchell et al., 2022), and model unlearning, which aims to remove the influence of specific data  
120 points (Bourtoule et al., 2021; Jang et al., 2023). Another line of work involves continually fine-  
121 tuning a pretrained model across a sequence of downstream tasks (Jin et al., 2021). Of particular  
122 relevance to our work is continual pretraining, where the pretraining process itself is extended to  
123 incorporate new data. This has been explored in domain-adaptive pretraining, in which models  
124 are sequentially trained on datasets from distinct, specialized domains (Gururangan et al., 2020;  
125 Chalkidis et al., 2020). However, these domains are often narrow in scope, and the datasets involved  
126 remain relatively small compared to those used in general pretraining. A notable exception is the  
127 work of Gupta et al. (2023), which studied the dynamics of training a large model on two datasets  
128 in sequence. Nevertheless, practical applications often require methods that scale to much longer  
129 sequences of datasets.

130 **Protein Language Models** Research in natural language processing (NLP) has recently been  
131 adapted to biology by treating the amino-acid sequence of proteins as a form of language. This  
132 perspective has led to the development of protein language models (pLMs), biologically inspired  
133 analogues of NLP models. For example, the autoregressive ProGen2 (Nijkamp et al., 2023) is  
134 based on GPT-2 (Radford et al.), while the masked ESM (Rives et al., 2021; Lin et al., 2023) and  
135 AMPLIFY (Fournier et al., 2024) draw inspiration from BERT (Devlin et al., 2019). Trained on  
136 large, diverse, and ever-growing protein sequence databases (Suzek et al., 2015; Jumper et al., 2021;  
137 Richardson et al., 2023), these models aim to capture evolutionary relationships and discover the  
138 underlying principles that govern protein structure and function. This approach has made pLMs an  
139 essential tool in computational biology for a wide range of applications such as mutational effect  
140 prediction, protein structure modeling, and de novo protein design (Hayes et al., 2025).

141 To evaluate the capabilities of protein language models, the community relies on several specialized  
142 benchmarks targeting different aspects of protein understanding. For protein folding, the Critical  
143 Assessment of protein Structure Prediction (CASP) is a biannual challenge that tests a model’s  
144 ability to predict 3D structures from amino acid sequences (J et al., 2018). For protein engineering  
145 and fitness prediction, the ProteinGym benchmark assesses how accurately models can predict the  
146 functional effects of mutations (Notin et al., 2023). In addition, broader multi-task benchmarks like  
147 TAPE (Rao et al., 2019) and PEER (Xu et al., 2022) evaluate model performance across a wide  
148 range of tasks, including remote homology detection and secondary structure prediction. In this  
149 work, we focus specifically on protein engineering and fitness prediction, given its crucial role in  
150 the drug discovery pipeline.

### 151 3 COPEP BENCHMARK

152  
153 To bridge the gap between continual learning research and its practical application, we introduce  
154 CoPeP, the Continual Pretraining for Protein Language Models benchmark. Built from successive  
155 UniProt releases, CoPeP reflects the challenge of keeping models updated with rapidly evolving  
156 biological data. It serves as a complex and large-scale testbed for continual learning methods, with  
157 significant implications for protein modeling and drug discovery.

#### 158 3.1 DATASET

159  
160 The CoPeP benchmark is constructed from the UniRef100 database (Suzek et al., 2015), which  
161 aggregates and clusters protein sequences curated by the UniProt Knowledgebase (The UniProt

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

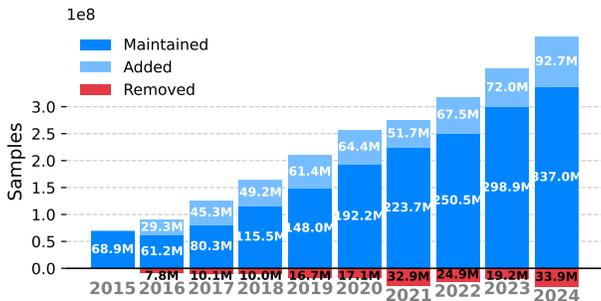


Figure 2: The number of proteins newly selected for and those no longer present in UniRef100 for each year in the benchmark. Despite millions of proteins being removed each year, the size of the dataset still grows as even more proteins are added.

Consortium, 2025), and constitutes the bulk of the training data for several pLMs (Rives et al., 2021; Lin et al., 2022; Fournier et al., 2024; Nijkamp et al., 2023). UniProt is updated multiple times each year, with millions of sequences added, removed, or replaced to reflect new biological knowledge and improved data quality. This evolving nature makes it an ideal foundation for evaluating continual pretraining.

For CoPeP, we select 10 consecutive yearly UniRef100 releases from 2015 to 2024, each corresponding to one task in the benchmark (the specific release and dates are listed in Table 4). These releases span hundreds of millions of protein sequences, with the dataset size increasing substantially year over year (Figure 2). Importantly, proteins may appear, disappear, or persist across releases: new sequences are introduced as biological discoveries accumulate, while others are removed if later deemed redundant or incorrect. Moreover, the dataset size does not grow linearly over time as, each year, an increasing number of new samples is added to the dataset.

Each sample is associated with an identifier and a protein sequence. However, the same identifier can map to multiple sequences, and vice-versa. To ensure consistency, we remove duplicate entries where both identifier and sequence are exact matches. Each dataset thus represents a faithful snapshot of the biological knowledge available at that time, capturing both growth in coverage and changes in curation practices. Together, these sequential datasets define the CoPeP training stream, providing a realistic setting to investigate how continual learning methods cope with evolving large-scale corpora.

### 3.2 STREAMING PROTOCOL

In traditional continual learning setups, training proceeds over a sequence  $\mathcal{D}_1, \dots, \mathcal{D}_n$  of  $n$  tasks, where each dataset  $\mathcal{D}_i = \{x_j\}_{j=1}^{m_i}$  is drawn from a task-specific data distribution  $x \sim \mathcal{P}_i$ . The challenge typically arises from distribution shifts between tasks, i.e.,  $\mathcal{P}_i \neq \mathcal{P}_{i+1}$ , which force the model to balance stability (retaining knowledge of earlier tasks) with plasticity (adapting to new tasks).

In CoPeP, the structure is slightly different. We also define a sequence  $\mathcal{D}_1, \dots, \mathcal{D}_n$  of  $n$  tasks, where each  $\mathcal{D}_i$  corresponds to the UniRef release from year  $i$ . However, in our case, these datasets are noisy snapshots of a common (unknown) underlying distribution  $\mathcal{P}^*$ . Importantly, the noise is systematic rather than random, as the protein datasets evolve over in a way reflecting community knowledge and interest. However, it is unknown how representative  $\mathcal{D}_i$  is of  $\mathcal{P}^*$ , with the challenge that yearly increments of the dataset do not correlate with improvements of  $\mathcal{P}_i$  w.r.t.  $\mathcal{P}^*$  (Fournier et al., 2024; Spinner et al., 2025).

Another difference between our setup and previous continual learning setups is that CoPeP does not forbid access to past data. Rather, at for task  $i$ , the learner may leverage the union of all observed datasets  $\mathcal{U}_i = \bigcup_{j=1}^i \mathcal{D}_j$ . This makes it possible to exploit temporal meta-information about the samples, such as the *multiplicity*  $c(x)$  of a sample, which counts how many consecutive years a protein has persisted in UniRef, i.e.,  $c(x) = \sum_{i=1}^k \mathbb{I}_{\mathcal{D}_i}(x)$ . Such information provides a signal

216 of sequence reliability, distinguishing consistently validated proteins from those that appear only  
217 transiently.

218 By structuring the problem this way, CoPeP reflects the practical challenges of maintaining large-  
219 scale models under real-world data evolution, while retaining the core challenges of continual learn-  
220 ing paradigms.

### 222 3.3 EVALUATION

224 Unlike traditional continual learning setups, because the underlying distribution that we are trying  
225 to learn is the same across all tasks, we are not concerned with metrics such as forgetting or transfer.  
226 Instead, at each evaluation timestep, we only measure the performance of the model on our suite of  
227 evaluation tasks at that specific timestep.

228 **Validation Set** We use the UniRef validation set introduced in Fournier et al. (2024) as part of  
229 our evaluations. These sequences were curated to be high-quality, complete proteins with strong  
230 experimental evidence for their existence. We deduplicated all of our training data against this  
231 validation set at the 90% sequence identity level using MMSeqs2 (Steinegger & Söding, 2017;  
232 Kallenborn et al., 2025) to ensure that the proteins in this validation set are not seen by the models  
233 at training time. The UniRef dataset contains proteins from all three domains of the phylogenetic  
234 tree of life (Bacteria, Archaea and Eukarya). Thus, performance on this set is an indicator of how  
235 well the model is able to reconstruct a broad range of proteins. We track both validation perplexity  
236 and accuracy.

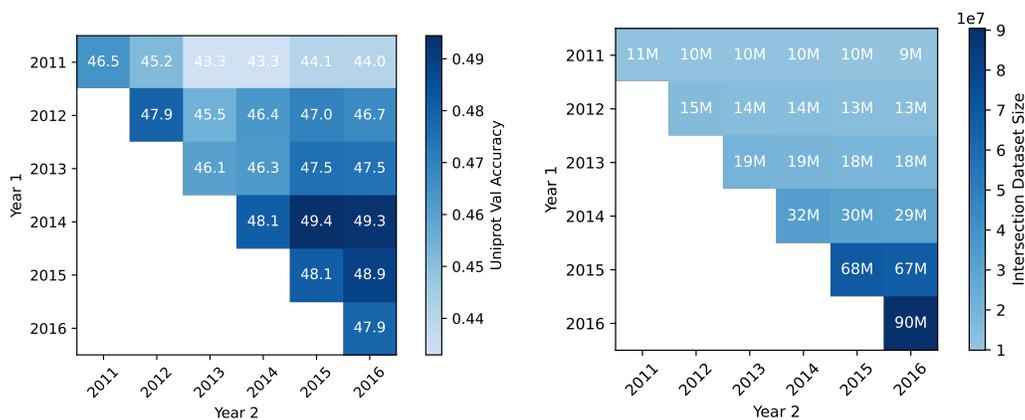
237 **ProteinGym** ProteinGym (Notin et al., 2023) is a broad benchmark designed for protein design and  
238 fitness prediction. It contains millions of mutated sequences from 217 deep mutational scanning  
239 assays across different taxa (humans, other eukaryotes, prokaryotes, and viruses). For each original  
240 sequence, the model ranks the mutations of that sequence by how likely they are, and this ranking  
241 is compared to ground truth rankings generated from experimental data and clinical annotations,  
242 computing the Spearman’s  $\rho$  between the two rankings. Across the set of assays, this results in more  
243 than 217 Spearman rank coefficients which we aggregate and report.

244 **PEER** PEER (Xu et al., 2022) is a multi-task benchmark for protein understanding. It contains 17  
245 tasks spanning various aspects of protein function prediction, protein localization prediction, protein  
246 structure prediction, protein-protein interaction prediction and protein-ligand interaction prediction.  
247 Because our model only operates on protein sequences, we exclude the protein-ligand interaction  
248 tasks, and we also exclude the ProteinNet based contact prediction task due to computational con-  
249 straints. This leaves us with 14 tasks from the original benchmark. Each task has its own evaluation  
250 metric, which we follow and report. We also report the average win rate of each model across all  
251 tasks compared to the other baselines, i.e. for each year, what percentage of the other models it  
252 outperforms.

253 **DGEB** DGEB (West-Roberts et al., 2024) is another multi-task benchmark designed to evaluate  
254 protein language models using diverse sequences across the tree of life and diverse tasks that capture  
255 different aspects of biological function. The benchmark contains 18 tasks, of which 16 use the amino  
256 acid modality. These tasks include classification, BiGene mining, Evolutionary Distance Similarity  
257 (EDS), pair classification, clustering, and retrieval. Similar to PEER, each task has its own evaluation  
258 metric, which we follow and report, along with the win rate for each method across all tasks.

### 260 3.4 BASE EXPERIMENTAL SETUP

261 The base model for all of our experiments is the AMPLIFY-120M (Fournier et al., 2024). It is  
262 an encoder model based on the BERT transformer (Devlin et al., 2019). For each task, we train  
263 for 100k steps using the AdamW optimizer (Loshchilov & Hutter, 2018) with weight decay set to  
264 .01 and an effective batch size of 4096. Fournier et al. (2024) use a cosine learning rate decay  
265 schedule, however, given the difficulty of rewarming up the learning rate after decay in continual  
266 learning (Gupta et al., 2023), we opt to use the warmup-stable-decay schedule (Hu et al., 2024; Li  
267 et al., 2025) which is more conducive to continuous training. For the first task, we linearly warm  
268 up the learning rate in the first 10k steps to .0005. At the end of each task, we linearly decay the  
269 learning rate to 0 over the final 10k steps of the task. When restarting training for the next task,  
we reset to the pre-decay checkpoint (i.e. the checkpoint right before the learning rate decay at 90k



(a) Validation accuracy on the UniProt validation dataset for the filtered datasets. The diagonals are unfiltered yearly releases, while each square in the top half shows the accuracy when the model was trained on the intersection of the data in the two years.

(b) Dataset sizes for filtered data experiments. Each square shows the number of protein sequences used in the training of the models in Figure 3a.

Figure 3: We train models on datasets that are the intersection of two yearly releases. Despite this filtering process creating smaller datasets, the validation accuracy actually improves for most years.

steps into training on the task). Thus, when starting on the sixth task in the benchmark, even though we have done 600k steps of training, the checkpoint we start with has only done 540k gradient steps.

### 3.5 USING TEMPORAL META-INFORMATION

As a preliminary experiment to validate the usefulness of the temporal meta-information described in Section 3.1, we test the hypothesis that data that stays in the UniProt database for longer is of higher quality and leads to better models. For this study, we take the releases for the first two years of our benchmark and releases from the four years prior to our benchmark (i.e. 2011-2016). For each pair of years, we only train on protein sequences that are in the intersection of both releases. Each model is trained for 100k steps according to the procedure outlined in Section 3.4. There are two hypothetical competing effects that this filtering could have: the smaller dataset size could lead to a decrease in performance or the potential increase in quality of data could lead to an increase in performance. We see the results in Figure 3. The performance of models trained on the unfiltered versions of the dataset are along the diagonal. From 2013 onwards, there is an increase in performance going from the unfiltered to filtered version of the datasets, even though the filtered datasets are smaller, implying that the benefit of the higher data quality wins out. For the first two years, since the datasets are already fairly small to begin with, filtering to an even smaller dataset has a deleterious effect. Across all years, however, we see there is an eventual increase in performance as you filter across a longer timespan. Furthermore, the best performance across all datasets comes from the intersection of 2014 and 2015 data, even though that dataset is a fraction of the size of the 2015 dataset, clearly showing the value of using temporal meta-information about the proteins.

## 4 METHODS

As shown in Section 3.5, the curation of data through the yearly updates of the UniProt Knowledgebase affects the prediction accuracy of the trained model. We now perform a large-scale study of 6 different methods to continually pretrain the AMPLIFY-120M model, that takes into account this temporality information. We focus on a set of representative methods spanning across 3 groups: continual learning, plasticity-focused and unlearning methods, with 2 algorithms for each group. Finally, we compare these methods with individual models trained on each yearly release separately, as is current standard practice (Fournier et al., 2024; Hayes et al., 2025) and with a model trained jointly on all data from 2015-2024.

#### 4.1 CONTINUAL LEARNING

**Sequential Training** This method is the simple baseline of training on each dataset in sequence, without any additional interventions or regularization. There are no additional hyperparameters for this method.

**Temporally Weighted Replay** Experience Replay (Rolnick et al., 2019; Abbes et al., 2025) is a commonly used technique in continual learning where a small subset of data from previous tasks is saved and rehearsed by the model while training on future tasks to prevent catastrophic forgetting. Given we can access all previous datasets and based on the results in Section 3.5, we use a modified version of this idea where we do not limit ourselves to a fixed size replay buffer. Instead, we continue sampling all samples according to how many previous datasets they appeared in. Let  $S = \{\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_{t-1}\}$  be the sequence of datasets up until the current task, and let  $U = \bigcup_{i=1}^{t-1} \mathcal{D}_i$  be their union. For any example  $x \in U$ , let its multiplicity be  $c(x) = \sum_{i=1}^k \mathbb{I}_{\mathcal{D}_i}(x)$ , where  $\mathbb{I}_{\mathcal{D}_i}(x)$  is the indicator function. The probability of sampling an example  $x$  from  $U$  is proportional to its multiplicity and is given by:  $P(x) = \frac{c(x)}{\sum_{y \in U} c(y)} = \frac{\sum_{i=1}^k \mathbb{I}_{\mathcal{D}_i}(x)}{\sum_{i=1}^k |\mathcal{D}_i|}$ . The total loss is given by  $(1 - \lambda_{replay})\mathcal{L}_{ce}(b_i) + \lambda_{replay}\mathcal{L}_{ce}(b_{replay})$ , where  $b_i$  is a batch sampled from the novel protein sequences added in task  $i$ ,  $b_{replay}$  a batch of samples from previous tasks sampled according to their multiplicity, and  $\lambda_{replay}$  weights the importance of current task compared to previous tasks.

#### 4.2 PLASTICITY

Loss of plasticity is a phenomenon in continual learning where as the model trains, it becomes less able to adapt to changes in data distributions. The plasticity preserving methods we use in our experiments are agnostic to the past data distributions and do not use any extra data.

**Shrink and Perturb** Shrink and Perturb (Ash & Adams, 2020) involves periodically shrinking and then adding noise to the weights of a neural network as a means of restoring plasticity to the network. In our experiments, at the start of every task, we set the weights as  $\theta_t = \lambda_{shrink}\theta_{t-1} + \lambda_{noise}p$ , where  $p$  are random weights drawn from the initialization distribution of the network.

**Hare and Tortoise** Hare and Tortoise (Lee et al., 2024) maintains two sets of network weights, slow and fast. The slow weights are an exponential moving average of the fast weights, i.e. at every step the slow weights are set to  $\theta_{slow} = \lambda_{ht.mom}\theta_{slow} + (1 - \lambda_{ht.mom})\theta_{fast}$ . Periodically, the fast weights are reset to the slow weights according to  $\lambda_{reset.freq}$ .

#### 4.3 UNLEARNING

Unlearning involves actively trying to remove knowledge about specific samples from the network. In our experiments, the forget set for task  $t$ ,  $\mathcal{F}_t$  is defined as the set of examples present in task  $t - 1$  but not in task  $t$ . With each step, we sample one batch from the current task  $b_i \sim \mathcal{D}_i$  and one batch from the forget set  $b_{forget} \sim \mathcal{F}_t$ .

**Gradient Ascent** Gradient ascent (Golatkar et al., 2020) attempts to unlearn knowledge by performing a gradient ascent step on data from the forget set. To prevent divergence, it also performs a descent step on data that is to be retained. This is implemented as optimizing the following loss:  $\mathcal{L}_{ce}(b_i) - \lambda_{asc}\mathcal{L}_{ce}(b_{forget})$  where  $\mathcal{L}_{ce}$  is the standard cross entropy loss used in training.

**Random Labels** Random labels (Golatkar et al., 2020) tries removing the knowledge in the forget set by sampling the targets of the forget set from the uniform distribution and performing gradient steps. The loss for the forget set is weighted by  $\lambda_{rand}$ .

#### 4.4 DESCRIPTION OF HYPERPARAMETER SEARCH AND OTHER EXPERIMENTAL DETAILS

For each method, we use the same base hyperparameters (e.g. learning rate, weight decay, batch size), and search over the method specific hyperparameters. Given the fact that several of these methods have not been used on such a scale, there does not exist much guidance in the literature on suitable ranges for many of these hyperparameters. We instead use an iterative, pruning based approach to our hyperparameter search to try a wide range for each method, and quickly prune suboptimal configurations. For each method, we evaluate 8 random configurations at 50k steps of

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

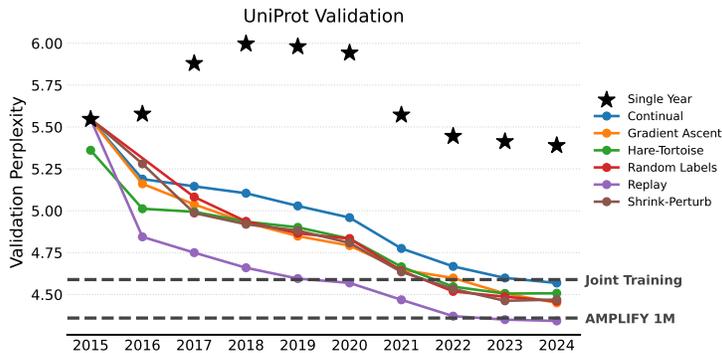


Figure 4: Validation perplexity on the UniProt validation set described in Section 3.3. All continual methods in our study beat the naive continual baseline and single year baseline.

total training. We then seed a Bayesian sampler with the results of those trials and sample 8 more configurations that are also evaluated at 50k steps of total training. The best 4 configurations from the 16 total trials are trained for another 150k steps, at which the best configuration is selected and trained for the remaining tasks in the benchmark. No method other than Hare and Tortoise deviates from the standard training for the first task in the benchmark, so every method (except for Hare and Tortoise) is started on task 2 from the pre-decay checkpoint of task 1. Hare and Tortoise is started from scratch on task 1. We use validation loss as the selection criterion in the search.

## 5 RESULTS

### 5.1 UNIPROT VALIDATION SET

In Figure 4, we show the performance of our models on the UniProt validation set. We notice several trends with our results. First, we see that performance generally seems to improve over time for the continual baselines. While this may seem trivial, it validates taking a continual approach to the problem. The steady improvement shows that continual training does not saturate the network or prime the network too heavily so that it cannot learn from future data. Furthermore, there is a big gap between the single year baseline and the continual baselines. This is partly because the continual models trained for longer, but this establishes that training from a continual checkpoint is effective. With each release, if the choice is to train for a certain number of steps from scratch or from a continual checkpoint, the continual checkpoint is a much more effective starting point. Surprisingly, every continual method also outperforms the joint training baseline that was trained on all data from 2015-2024 for the same number of steps as the continual models. This is likely because the joint training baseline also learns from data that was removed, while the continual models only learn from data is present in the current release.

We should also note that several models start reaching the performance level of AMPLIFY 1M (the base model trained for 1 million steps according to Fournier et al. (2024)) with considerably fewer steps and access to less data throughout training. In fact, the temporal replay method essentially matches the performance of AMPLIFY-1M at 8 tasks, which is the equivalent of 730k steps, with much of the training taking place with access to much less data, and overtakes it after the 9th task.

Finally, comparing the methods amongst each other, we see that every method offers better performance compared to the naive continual baseline and (other than the temporal replay baseline) relatively similar performance to each other. This is highly encouraging, as essentially none of these methods were developed for this specific setup, and yet they are all showing positive performance. Hare and Tortoise and Shrink and Perturb are both plasticity preserving methods, but to our knowledge have never been applied to a model or training scale of this size. Gradient Ascent and Random Labels have never been used with LLMs, but generally on more limited forget sets and not as a part of continual pretraining. The relative success of these methods shows that all of these approaches to continual learning (forgetting, plasticity, unlearning) have ideas to contribute in this setup.

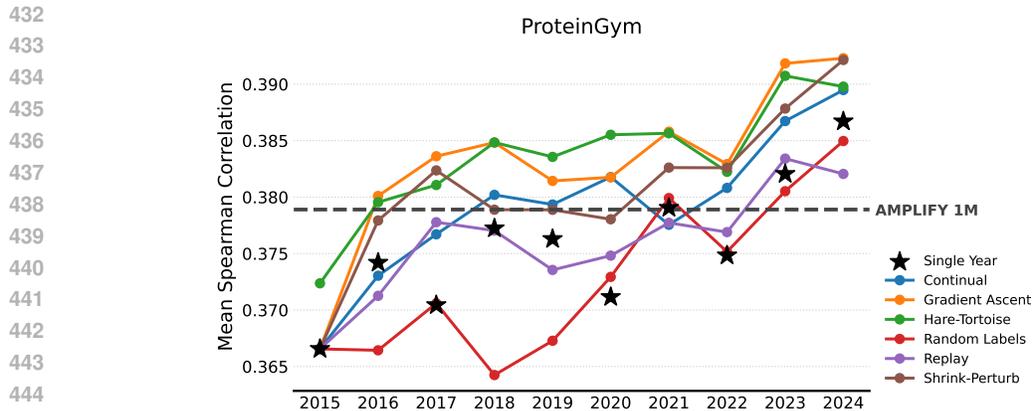


Figure 5: Results on the ProteinGym benchmark. Several continual methods outperform the naive continual baseline, as well as the AMPLIFY-1M baseline, a version of our model that was trained for 1 million steps.

## 5.2 PROTEINGYM

The trends for the ProteinGym evaluation (Figure 5) are slightly harder to define compared to the UniProt Validation set. Gradient Ascent and Shrink and Perturb both perform the best across all methods. All the continual and the single year baseline at 2024 outperforms AMPLIFY-1M, with 3 of the non-naive continual learning methods also outperforming naive continual learning. Surprisingly, the temporal replay method, which does the best on the UniProt validation set, does performs the worst on ProteinGym. We discuss the potential reasons for this in Section 5.4. There does seem to be a measure of early saturation for a while, as the improvement in performance for most of the methods seems to plateau after the first 3 tasks, but then we see a jump in performance at the end when the 2023 and 2024 data is introduced. This is likely because ProteinGym was introduced in 2023, and thus many of the proteins in the benchmark were likely added to the UniProt database in the last two years.

## 5.3 MULTITASK PROTEIN UNDERSTANDING BENCHMARKS

We show the a summary of the results on the multitask protein understanding benchmarks DGEB and PEER in Figures 6a and 6b respectively. The win rates are computed compared to every other checkpoint for all methods and years, across all tasks, and thus is a relatively stable metric. For more fine-grained results on these benchmarks, please see Appendices C and D.

For PEER 6a, none of the Single Year baselines perform significantly better than all the rest, implying that there is not necessarily a year that is better for downstream performance. Some Single Year results are worse than average, however, and it is notable that after training on 2022 data (one of the worst Single Year performances), the performance of the Gradient Ascent method essentially collapsed. This is also seen in the DGEB results 6b. The best performing model at any point was Shrink and Perturb at years 2020-2021, followed by Temporal Replay at year 2023.

For DGEB 6b, we see that the naive continual baseline does not perform well, with nearly all continual learning methods outperforming it (except for Gradient Ascent after training on 2022 data). Interestingly, the best model again is at year 2020, although in this case it is Random Labels that achieves the best performance. Temporal Replay again achieves strong performance at year 2023. Overall, the results for DGEB generally align with the PEER results, but both of these benchmarks seem to be not very aligned with the ProteinGym results.

## 5.4 TRADEOFFS

Drug discovery is a long process and requires many different capabilities with respect to proteins, including generation, property prediction, fitness prediction, and optimization. It is difficult to create

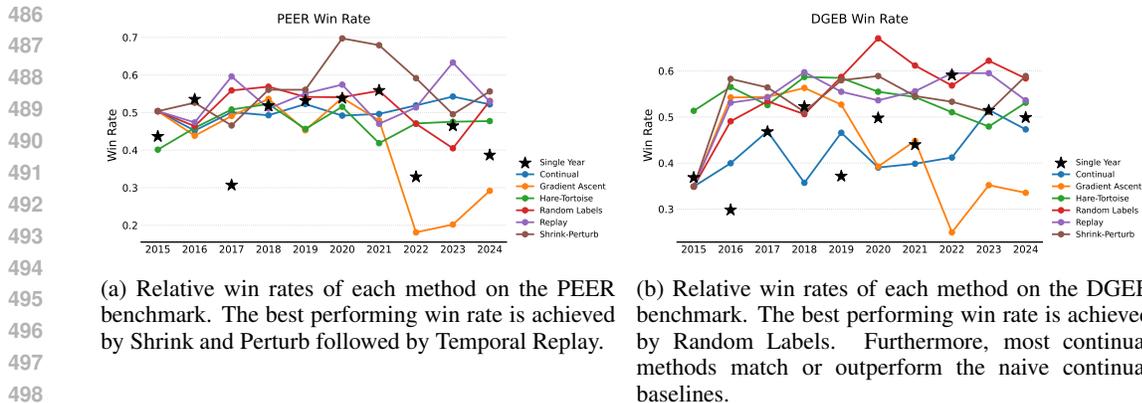


Figure 6: Results on the multitask protein understanding benchmarks.

an evaluation that is able to cover all of these capabilities. In this section we discuss the tradeoffs of the two evaluations we use in our benchmark.

Both the UniProt validation set and ProteinGym were curated from natural proteins, which means that models that do well on them would be helpful in creating therapeutic drugs, but not necessarily industrial proteins. The UniProt validation set was constructed to be as close as possible to the natural distribution of proteins, with the idea that nature is a good inductive bias. The proteins were selected to be from diverse and highly studied proteomes. A consequence of the latter point is that the highly studied proteomes were likely in the earlier UniRef100 releases which could explain the success of a method like temporal replay that upweights such samples. Because we deduplicated our training set against the validation set, however, the evaluation rewards models that do not overfit to specific samples or mutations that they see in the data, and instead learn to generalize to the larger patterns. On the other hand, ProteinGym rewards models that can properly evaluate the fitness of specific mutations in a protein. Given we did not deduplicate our training set against ProteinGym and that memorizing specific sequences could provide an advantage to the model, it is also possible that ProteinGym would reward models that overfit the data slightly.

## 6 DISCUSSION

We present CoPeP, a benchmark for continual pretraining of protein language models. The datasets used in our benchmark are curated from the regular releases of UniProt, and thus naturally evolve as the biologist community’s knowledge and interest evolve. CoPeP is regularly extensible as each new UniProt release becomes available, making it more difficult to saturate the benchmark. In our work, we show that several different approaches to continual learning and unlearning are able to improve on naive continual learning, and our benchmark is an opportunity for those communities to develop and test their methods on a realistic, large scale setting. Several of the methods we present are also fairly orthogonal to each other, and future work can investigate how to combine them to create a better method. Although not explored in our work, the closely related field of model editing could also potentially apply contribute to this problem.

Our work also explores the idea of using temporal meta-information about each sample to guide training. We use this information as both a filter and as a replay strategy, and show that both approaches improves performance. Future work should explore protein specific learning methods that can better leverage this temporal meta-information.

We hope that this benchmark can accelerate progress in protein language model learning. For large biomedical companies, it may be cost-feasible to simply retrain from scratch on large data, but having to do so takes time that can lengthen experiment cycles. Effective continual training could also enable academic labs to perform relevant and cost-effective research and further push the frontier of drug discovery.

## 7 REPRODUCIBILITY STATEMENT

The details of our model training and hyperparameter selection are provided in Sections 3.4, 4.4, F.1, and G. The details of our dataset curation are provided in Sections 3.1, 3.3, and H. Upon acceptance, we also intend to release the code, checkpoints, and datasets used to conduct all of our experiments.

## REFERENCES

- Istabrak Abbes, Gopeshh Subbaraj, Matthew Riemer, Nizar Islah, Benjamin Therien, Tsuguchika Tabaru, Hiroaki Kingetsu, Sarath Chandar, and Irina Rish. Revisiting Replay and Gradient Alignment for Continual Pre-Training of Large Language Models, August 2025.
- Jordan Ash and Ryan P Adams. On Warm-Starting Neural Network Training. In *Advances in Neural Information Processing Systems*, volume 33, pp. 3884–3894. Curran Associates, Inc., 2020.
- Jorg Bornschein, Alexandre Galashov, Ross Hemsley, Amal Rannen-Triki, Yutian Chen, Arslan Chaudhry, Xu Owen He, Arthur Douillard, Massimo Caccia, Qixuan Feng, et al. Nevis’ 22: A stream of 100 tasks sampled from 30 years of computer vision research. *Journal of Machine Learning Research*, 24(308):1–77, 2023.
- Lucas Bourtole, Varun Chandrasekaran, Christopher A. Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine Unlearning. In *2021 IEEE Symposium on Security and Privacy (SP)*, pp. 141–159, May 2021. doi: 10.1109/SP40001.2021.00019.
- Zhipeng Cai, Ozan Sener, and Vladlen Koltun. Online continual learning with natural distribution shifts: An empirical study with visual data. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 8281–8290, 2021.
- Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. LEGAL-BERT: The Muppets straight out of Law School. In Trevor Cohn, Yulan He, and Yang Liu (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 2898–2904, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.261.
- Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K. Dokania, Philip H. S. Torr, and Marc’ Aurelio Ranzato. On Tiny Episodic Memories in Continual Learning, June 2019.
- Peter J. A. Cock, Tiago Antao, Jeffrey T. Chang, Brad A. Chapman, Cymon J. Cox, Andrew Dalke, Iddo Friedberg, Thomas Hamelryck, Frank Kauff, Bartek Wilczynski, and Michiel J. L. de Hoon. Biopython: Freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*, 25(11):1422–1423, June 2009. ISSN 1367-4803. doi: 10.1093/bioinformatics/btp163.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423.
- Shibhansh Dohare, J Fernando Hernandez-Garcia, Qingfeng Lan, Parash Rahman, A Rupam Mahmood, and Richard S Sutton. Loss of plasticity in deep continual learning. *Nature*, 632(8026): 768–774, 2024.
- Quentin Fournier, Robert M. Vernon, Almer van der Sloot, Benjamin Schulz, Sarath Chandar, and Christopher James Langmead. Protein Language Models: Is Scaling Necessary?, September 2024.

- 594 Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal Sunshine of the Spotless Net:  
595 Selective Forgetting in Deep Networks. In *2020 IEEE/CVF Conference on Computer Vision and*  
596 *Pattern Recognition (CVPR)*, pp. 9301–9309, Seattle, WA, USA, June 2020. IEEE. ISBN 978-1-  
597 7281-7168-5. doi: 10.1109/CVPR42600.2020.00932.
- 598
- 599 Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empiri-  
600 cal investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint*  
601 *arXiv:1312.6211*, 2013.
- 602
- 603 Kshitij Gupta, Benjamin Thérien, Adam Ibrahim, Mats L. Richter, Quentin Anthony, Eugene  
604 Belilovsky, Irina Rish, and Timothée Lesort. Continual Pre-Training of Large Language Mod-  
605 els: How to (re)warm your model?, September 2023.
- 606
- 607 Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,  
608 and Noah A. Smith. Don’t Stop Pretraining: Adapt Language Models to Domains and Tasks.  
609 In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th*  
610 *Annual Meeting of the Association for Computational Linguistics*, pp. 8342–8360, Online, July  
611 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.740.
- 612
- 613 Thomas Hayes, Roshan Rao, Halil Akin, Nicholas J Sofroniew, Deniz Oktay, Zeming Lin, Robert  
614 Verkuil, Vincent Q Tran, Jonathan Deaton, Marius Wiggert, et al. Simulating 500 million years  
615 of evolution with a language model. *Science*, 387(6736):850–858, 2025.
- 616
- 617 Shengding Hu, Yuge Tu, Xu Han, Ganqu Cui, Chaoqun He, Weilin Zhao, Xiang Long, Zhi Zheng,  
618 Yewei Fang, Yuxiang Huang, Xinrong Zhang, Zhen Leng Thai, Chongyi Wang, Yuan Yao,  
619 Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chao Jia, Guoyang Zeng, Da-  
620 hai Li, Zhiyuan Liu, and Maosong Sun. MiniCPM: Unveiling the Potential of Small Language  
621 Models with Scalable Training Strategies. In *First Conference on Language Modeling*, August  
622 2024.
- 623
- 624 Moul J, Fidelis K, Kryshchak A, Schwede T, and Tramontano A. Critical assessment of methods  
625 of protein structure prediction (CASP)-Round XII. *Proteins*, 2018. doi: 10.1002/prot.25415.
- 626
- 627 Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and  
628 Minjoon Seo. Knowledge Unlearning for Mitigating Privacy Risks in Language Models. In  
629 Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual*  
630 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14389–  
631 14408, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/  
632 v1/2023.acl-long.805.
- 633
- 634 Xisen Jin, Bill Yuchen Lin, Mohammad Rostami, and Xiang Ren. Learn Continually, Gener-  
635 alize Rapidly: Lifelong Knowledge Accumulation for Few-shot Learning. In Marie-Francine  
636 Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Findings of the Associ-*  
637 *ation for Computational Linguistics: EMNLP 2021*, pp. 714–729, Punta Cana, Dominican Re-  
638 public, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.  
639 findings-emnlp.62.
- 640
- 641 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,  
642 Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland,  
643 Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-  
644 Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman,  
645 Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Se-  
646 bastian Bodenstern, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Push-  
647 meet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with AlphaFold.  
*Nature*, 596(7873):583–589, August 2021. ISSN 1476-4687. doi: 10.1038/s41586-021-03819-2.
- 648
- 649 Felix Kallenborn, Alejandro Chacon, Christian Hundt, Hassan Sirelkhatim, Kieran Didi, Sooyoung  
650 Cha, Christian Dallago, Milot Mirdita, Bertil Schmidt, and Martin Steinegger. GPU-accelerated  
651 homology search with MMseqs2. *Nat Methods*, pp. 1–4, September 2025. ISSN 1548-7105. doi:  
652 10.1038/s41592-025-02819-8.

- 648 Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. Continual Pre-  
649 training of Language Models. In *The Eleventh International Conference on Learning Representations*,  
650 September 2022.
- 651 James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A  
652 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcom-  
653 ing catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*,  
654 114(13):3521–3526, 2017.
- 655 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-  
656 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A  
657 benchmark of in-the-wild distribution shifts. In *International conference on machine learning*,  
658 pp. 5637–5664. PMLR, 2021.
- 659 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
660 2009.
- 661 Hojoon Lee, Hyeonseo Cho, Hyunseung Kim, Donghu Kim, Dugki Min, Jaegul Choo, and Clare  
662 Lyle. Slow and Steady Wins the Race: Maintaining Plasticity with Hare and Tortoise Networks,  
663 June 2024.
- 664 Yunshui Li, Yiyuan Ma, Shen Yan, Chaoyi Zhang, Jing Liu, Jianqiao Lu, Ziwen Xu, Mengzhao  
665 Chen, Minrui Wang, Shiyi Zhan, Jin Ma, Xunhao Lai, Yao Luo, Xingyan Bin, Hongbin Ren,  
666 Mingji Han, Wenhao Hao, Bairen Yi, LingJun Liu, Bole Ma, Xiaoying Jia, Zhou Xun, Liang  
667 Xiang, and Yonghui Wu. Model Merging in Pre-training of Large Language Models, May 2025.
- 668 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Allan dos Santos  
669 Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, and Alexander Rives. Language models  
670 of protein sequences at the scale of evolution enable accurate structure prediction, July 2022.
- 671 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,  
672 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom  
673 Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level  
674 protein structure with a language model. *Science*, 379(6637):1123–1130, March 2023. doi:  
675 10.1126/science.ade2574.
- 676 Zhiqiu Lin, Jia Shi, Deepak Pathak, and Deva Ramanan. The clear benchmark: Continual learn-  
677 ing on real-world imagery. In *Thirty-fifth conference on neural information processing systems*  
678 *datasets and benchmarks track (round 2)*, 2021.
- 679 Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *International Con-*  
680 *ference on Learning Representations*, September 2018.
- 681 Ali Madani, Ben Krause, Eric R Greene, Subu Subramanian, Benjamin P Mohr, James M Holton,  
682 Jose Luis Olmos Jr, Caiming Xiong, Zachary Z Sun, Richard Socher, et al. Large language  
683 models generate functional protein sequences across diverse families. *Nature biotechnology*, 41  
684 (8):1099–1106, 2023.
- 685 Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The  
686 sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165.  
687 Elsevier, 1989.
- 688 Leland McInnes, John Healy, Nathaniel Saul, and Lukas Grossberger. Umap: Uniform manifold  
689 approximation and projection. *The Journal of Open Source Software*, 3(29):861, 2018.
- 690 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual  
691 associations in GPT. In *Proceedings of the 36th International Conference on Neural Information*  
692 *Processing Systems, NIPS ’22*, pp. 17359–17372, Red Hook, NY, USA, November 2022. Curran  
693 Associates Inc. ISBN 978-1-7138-7108-8.
- 694 Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. Memory-  
695 Based Model Editing at Scale. In *Proceedings of the 39th International Conference on Machine*  
696 *Learning*, pp. 15817–15831. PMLR, June 2022.

- 702 Erik Nijkamp, Jeffrey A. Ruffolo, Eli N. Weinstein, Nikhil Naik, and Ali Madani. ProGen2: Ex-  
703 ploring the boundaries of protein language models. *Cell Systems*, 14(11):968–978.e3, November  
704 2023. ISSN 2405-4712. doi: 10.1016/j.cels.2023.10.002.
- 705  
706 Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood Van Niekerk, Steffanie Paul, Han Spinner, Nathan  
707 Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, et al. Proteingym: Large-scale benchmarks  
708 for protein fitness prediction and design. *Advances in Neural Information Processing Systems*, 36:  
709 64331–64379, 2023.
- 710 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language  
711 Models are Unsupervised Multitask Learners.
- 712  
713 Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Xi Chen, John Canny, Pieter Abbeel,  
714 and Yun S. Song. Evaluating protein transfer learning with TAPE. In *Proceedings of the 33rd*  
715 *International Conference on Neural Information Processing Systems*, number 869, pp. 9689–  
716 9701. Curran Associates Inc., Red Hook, NY, USA, December 2019.
- 717 Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl:  
718 Incremental classifier and representation learning. In *Proceedings of the IEEE conference on*  
719 *Computer Vision and Pattern Recognition*, pp. 2001–2010, 2017.
- 720 Lorna Richardson, Ben Allen, Germana Baldi, Martin Beracochea, Maxwell L Bileschi, Tony Bur-  
721 dett, Josephine Burgin, Juan Caballero-Pérez, Guy Cochrane, Lucy J Colwell, Tom Curtis, Alejandra  
722 Escobar-Zepeda, Tatiana A Gurbich, Varsha Kale, Anton Korobeynikov, Shriya Raj, Alexander  
723 B Rogers, Ekaterina Sakharova, Santiago Sanchez, Darren J Wilkinson, and Robert D Finn.  
724 MGnify: The microbiome sequence data analysis resource in 2023. *Nucleic Acids Res*, 51(D1):  
725 D753–D759, January 2023. ISSN 0305-1048. doi: 10.1093/nar/gkac1080.
- 726 Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo,  
727 Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function  
728 emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings*  
729 *of the National Academy of Sciences*, 118(15):e2016239118, April 2021. doi: 10.1073/pnas.  
730 2016239118.
- 731 David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experi-  
732 ence Replay for Continual Learning. In *Advances in Neural Information Processing Systems*,  
733 volume 32. Curran Associates, Inc., 2019.
- 734  
735 Aviv Spinner, Erika DeBenedictis, and Corey M. Hudson. Scaling and Data Saturation in Protein  
736 Language Models, July 2025.
- 737 Martin Steinegger and Johannes Söding. MMseqs2 enables sensitive protein sequence searching  
738 for the analysis of massive data sets. *Nat Biotechnol*, 35(11):1026–1028, November 2017. ISSN  
739 1546-1696. doi: 10.1038/nbt.3988.
- 740  
741 Baris E. Suzek, Yuqi Wang, Hongzhan Huang, Peter B. McGarvey, Cathy H. Wu, and the UniProt  
742 Consortium. UniRef clusters: A comprehensive and scalable alternative for improving sequence  
743 similarity searches. *Bioinformatics*, 31(6):926–932, March 2015. ISSN 1367-4803. doi: 10.1093/  
744 bioinformatics/btu739.
- 745  
746 The UniProt Consortium. UniProt: The Universal Protein Knowledgebase in 2025. *Nucleic Acids*  
747 *Res*, 53(D1):D609–D617, January 2025. ISSN 1362-4962. doi: 10.1093/nar/gkae1010.
- 748  
749 Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A Comprehensive Survey of Continual  
750 Learning: Theory, Method and Application. *IEEE Trans. Pattern Anal. Mach. Intell.*, 46(8):  
751 5362–5383, August 2024. ISSN 0162-8828. doi: 10.1109/TPAMI.2024.3367329.
- 752  
753 Jacob West-Roberts, Joshua Kravitz, Nishant Jha, Andre Cornman, and Yunha Hwang. Diverse  
754 Genomic Embedding Benchmark for functional evaluation across the tree of life, July 2024.
- 755  
756 Minghao Xu, Zuobai Zhang, Jiarui Lu, Zhaocheng Zhu, Yangtian Zhang, Chang Ma, Runcheng  
757 Liu, and Jian Tang. PEER: A Comprehensive and Multi-Task Benchmark for Protein Sequence  
758 Understanding. In *Thirty-Sixth Conference on Neural Information Processing Systems Datasets*  
759 *and Benchmarks Track*, June 2022.

756 Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen,  
757 and Ningyu Zhang. Editing Large Language Models: Problems, Methods, and Opportunities.  
758 In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on*  
759 *Empirical Methods in Natural Language Processing*, pp. 10222–10240, Singapore, December  
760 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.632.

761 Çağatay Yıldız, Nishaanth Kanna Ravichandran, Nitin Sharma, Matthias Bethge, and Beyza Er-  
762 mis. Investigating Continual Pretraining in Large Language Models: Insights and Implications.  
763 *Transactions on Machine Learning Research*, February 2025. ISSN 2835-8856.

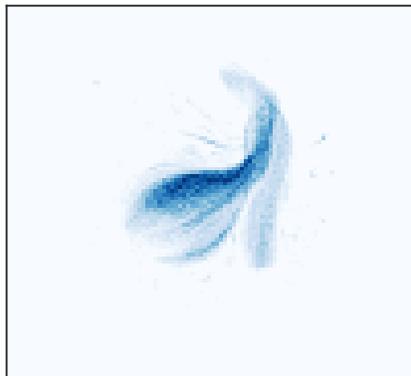
764 Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence.  
765 In *International conference on machine learning*, pp. 3987–3995. PMLR, 2017.

766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

## 810 A EVOLUTION OF DATA

### 811 A.1 EMBEDDING VISUALIZATION

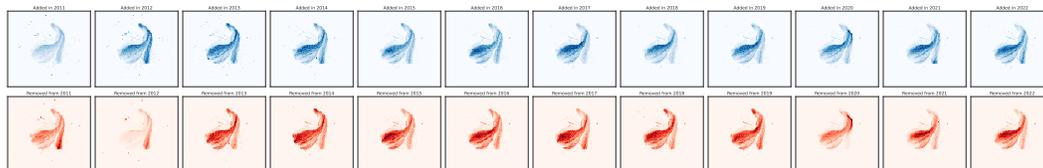
812 In this experiment, we analyze how protein sequence datasets evolve over years by visualizing their  
 813 structure in embedding space. We use representations from AMPLIFY (trained with 1 million steps)  
 814 and apply UMAP (McInnes et al., 2018) to project high-dimensional protein embeddings into two  
 815 dimensions. This enables us to observe broad patterns in the data and how they change across  
 816 consecutive UniRef100 releases.



817 Figure 7: UMAP projection of protein embeddings from all UniRef100 releases (AMPLIFY 1M  
 818 representations). The plot shows a stable global structure with a dense core and branches, indicating  
 819 natural groupings of proteins.

820 Figure 7 shows the embedding of the full dataset i.e., sequences from all years. The plot reveals a  
 821 global structure with a dense central region and branches, suggesting natural groupings of proteins.  
 822 Differences in density highlight areas where certain types of sequences are more common.

823 Each UniRef100 release both adds and removes sequences, reflecting the expansion of biological  
 824 knowledge and ongoing curation. To illustrate these dynamics, Figure 8 compares additions (blue,  
 825 top row) with removals (red, bottom row) per year. Overall, while the global structure of protein  
 826 embeddings is stable, Figure 8 indicates local shifts such as density increases and cluster expansion.  
 827 This underscores why continual learning is critical for protein language models. Instead of treating  
 828 each release as an isolated datasets, continual methods can exploit temporal information to adapt to  
 829 new proteins as well as retain knowledge.



830 Figure 8: Yearly dynamics of UniRef100 embeddings. Top row (blue): proteins added in each year;  
 831 bottom row (red): proteins removed. While the global organization of protein embeddings is stable,  
 832 the local shifts such as density increases and cluster expansion are indicate yearly shift in underlying  
 833 distribution.

### 834 A.2 MODEL EMBEDDING SHIFTS

835 In Figure 9, we visualize how different continual learning methods structure their protein embed-  
 836 dings. For each method, we select the final checkpoint, and extract embeddings for a representative  
 837 subset of proteins that were added and removed across all years. We then apply UMAP to project  
 838 these embeddings into two dimensions. The color shading indicates the density of proteins from

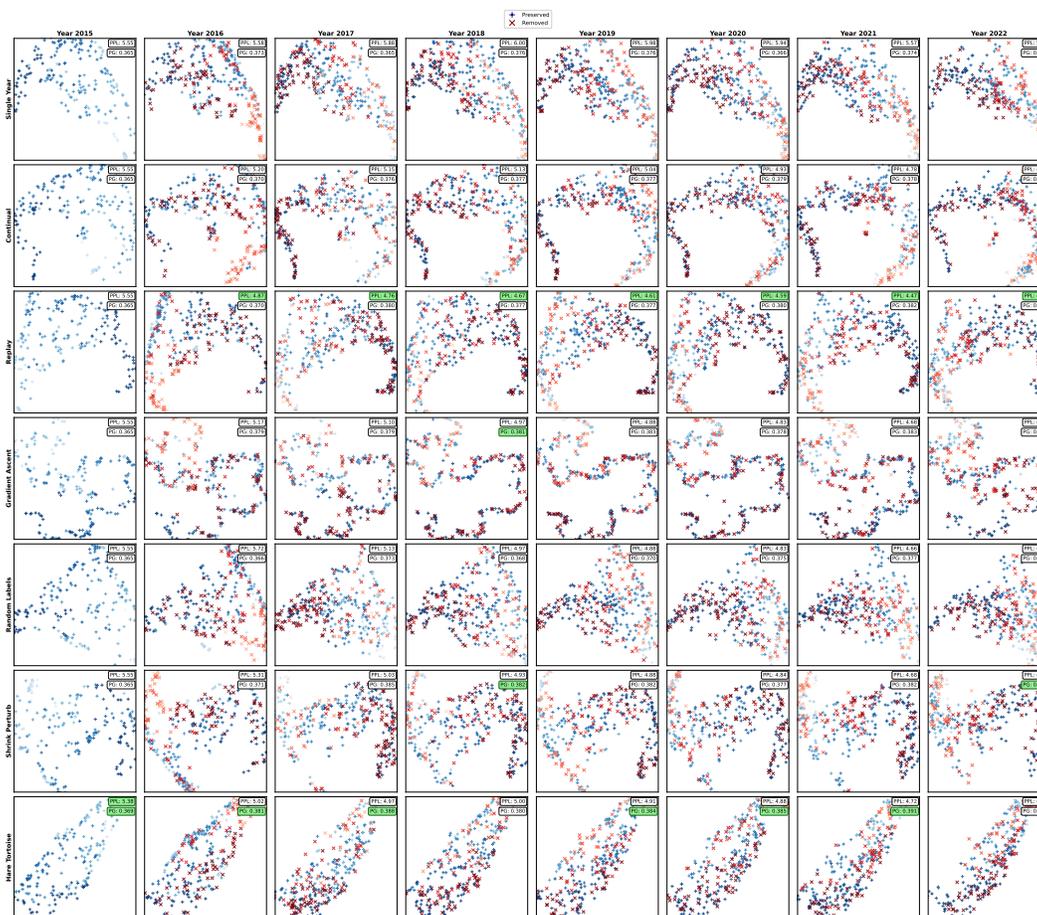


Figure 9: UMAP visualization of protein embeddings across continual learning methods (2015–2022) with performance metrics (PPL: validation perplexity; PG: ProteinGym mean Spearman correlation). Colors shades are based on single year clusters. Darker shades indicate denser parts. Same color shades are used in other rows/methods to indicate those centroid/representative proteins locations.

the single year clusters, with darker shades representing denser regions. Interestingly, most methods maintain a similar global structure except for Gradient Ascent, which shows a more rope-like structure.

### A.3 DATA STATISTICS

In this section, we analyze how various protein sequence statistics evolve over years. We compute several statistics for each protein sequence in the UniRef100 releases from 2015 to 2024 using the Biopython library (Cock et al., 2009). Specifically, we create two sample QQ-plots comparing the distribution of each statistic between each year and the reference year 2015. The statistics we analyze include:

- Aromaticity
- Charge at pH 7
- Instability Index
- Isoelectric Point
- Molar Extinction Coefficient (oxidized and reduced)
- Protein Length
- Longest Repeat Ratio

Overall, we see that most statistics show relatively stable distributions across years, for at least some of the statistics (molar extinction coefficient, protein length, charge at pH 7) there are outliers at

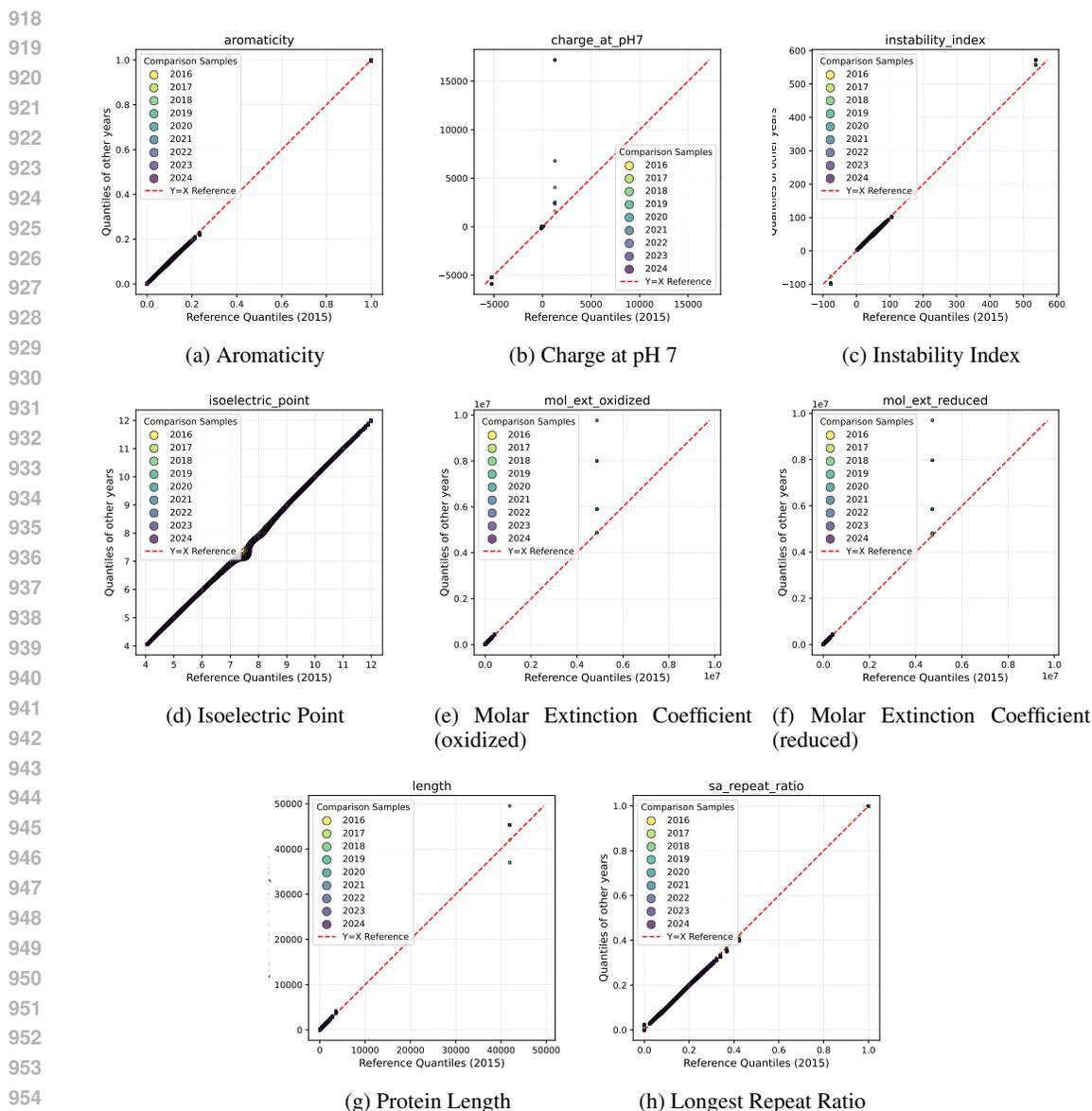


Figure 10: Two sample QQ-plots of various protein sequence statistics across different years. Each line represents a comparison between a year and the reference year 2015.

the extremes of the distributions for several years. Furthermore, there is a noticeable shift in the distribution of the isoelectric point over the years, with later years showing a higher density of proteins with isoelectric points around 7-8.

## B FINE GRAINED RESULTS ON PROTEINGYM

We also visualized two variants of the ProteinGym evaluation, similar to Figure 5: the best performance achieved in each year of training and the fine-grained trajectory of performance across all steps.

In Figure 11, we observe that Hare Tortoise consistently delivers the strongest results, with Gradient Ascent and Shrink Perturb close behind. All three methods perform better than the AMPLIFY 1M baseline across nearly all year releases, while continual learning and replay result in modest gains.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

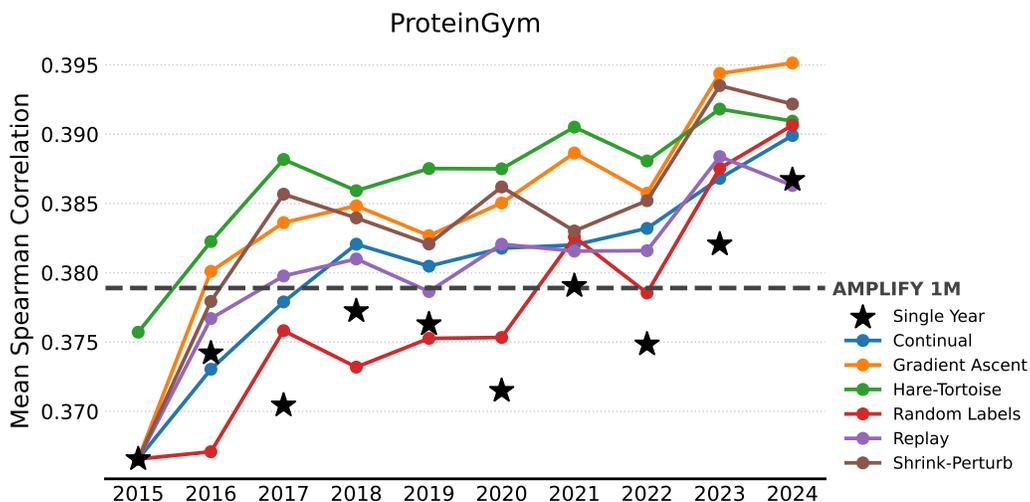


Figure 11: Best mean Spearman correlation for continual training on ProteinGym. Hare Tortoise achieves the best performance across nearly all years, with Gradient Ascent and Shrink Perturb close behind. These methods consistently perform better than AMPLIFY 1M.

Random Labels again shows improvement relative to the naive Single Year baseline, but it does not reach the same level as the other methods.

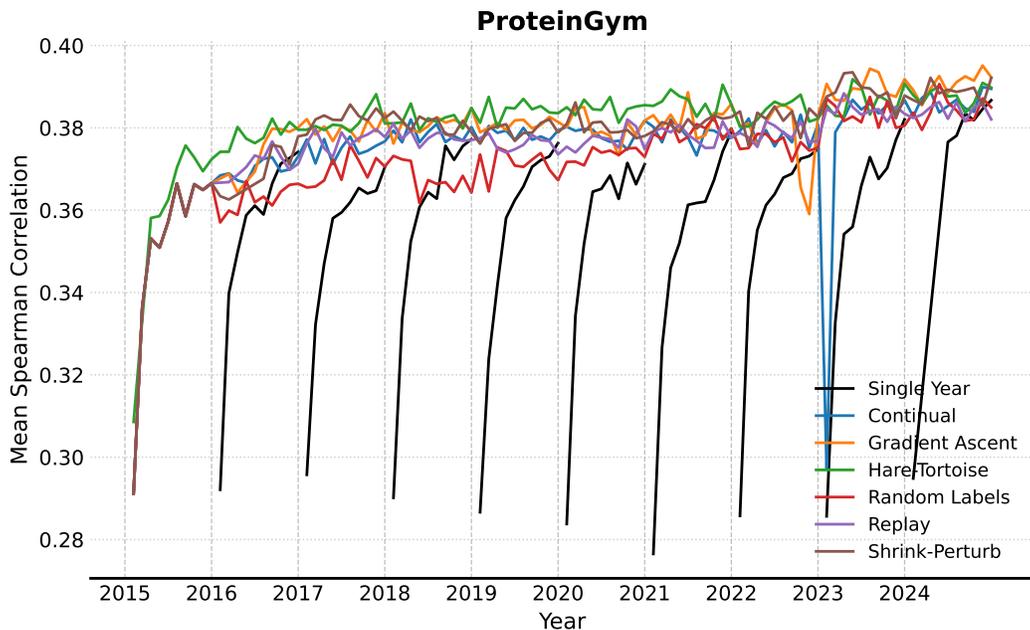


Figure 12: Mean Spearman correlation on ProteinGym across training steps for continual training on ProteinGym. While the naive Single Year baseline resets each year and lags substantially, continual learning methods such as Hare Tortoise, Gradient Ascent, and Shrink Perturb maintain strong performance throughout training and consistently perform better than AMPLIFY 1M.

Figure 12 highlight the shortcomings of Single Year most clearly as the model start from scratch. By contrast, Hare Tortoise, Gradient Ascent, and Shrink Perturb maintain strong performance throughout training, suggesting that these methods provide more stable and reliable learning dynamics.

Apart from these results, we also provide the boxplots of Spearman correlations across methods in Figure 13. In all cases, Hare Tortoise, Gradient Ascent, and Shrink Perturb consistently cluster above the AMPLIFY 1M baseline, with relatively tight distributions indicating robust improvements. Replay and Continual show more variance, often overlapping with the baseline but generally outperforming the naive Single Year approach.

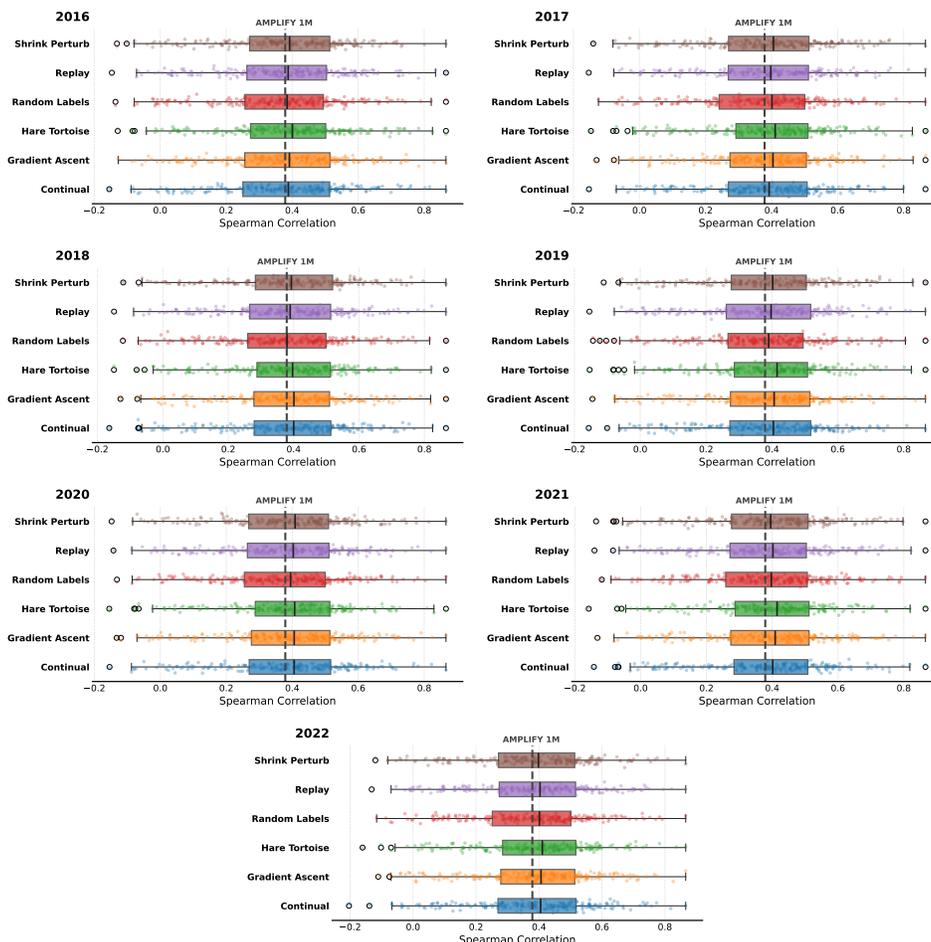


Figure 13: Distribution of Spearman correlations for each method from 2016–2022 on ProteinGym benchmarks. Hare Tortoise, Gradient Ascent, and Shrink Perturb consistently center above the AMPLIFY 1M baseline with tight variance, indicating strong and stable performance.

Table 1 summarizes the AUC performance of different methods on ProteinGym. We again observe that Hare Tortoise consistently achieves the best or tied-best results across nearly all years, with Gradient Ascent and Shrink Perturb closely following. These findings align with our observations in Figure 5. By contrast, continual learning and replay provide moderate gains over the naive Single Year baseline.

### C FINE GRAINED RESULTS ON PEER

In Figure 14, we show the full results for each task on the PEER benchmark.

### D FINE GRAINED RESULTS ON DGEB

In Figure 15, we show the full results for each task on the DGEB

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

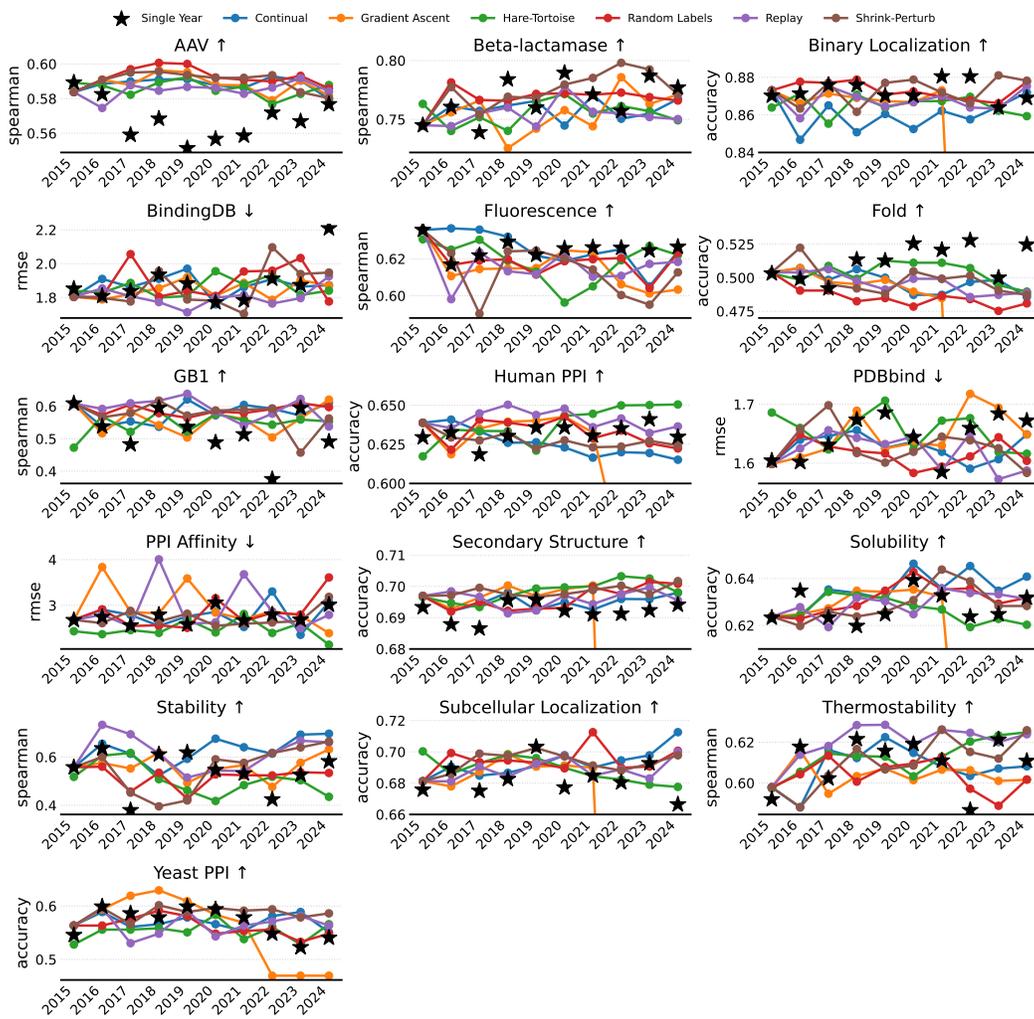


Figure 14: Full results on the PEER benchmark.

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

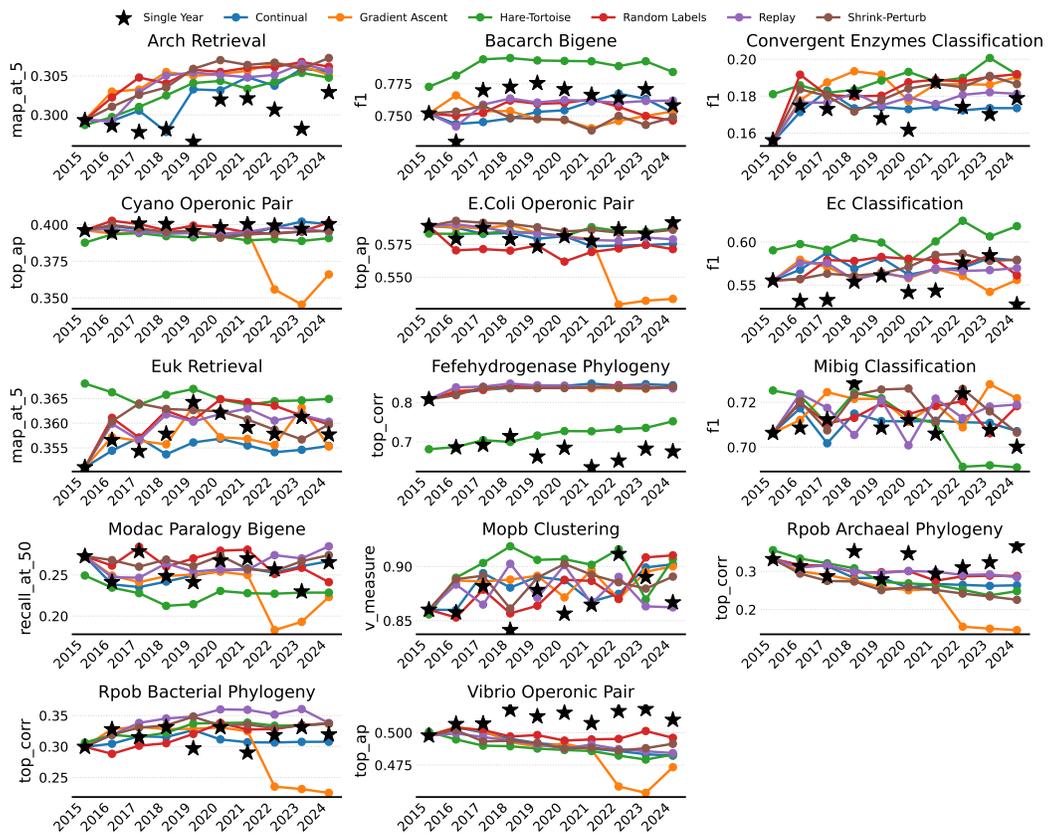


Figure 15: Full results on the DGEb benchmark.

	Single Year	Continual	Gradient Ascent	Hare Tortoise	Random Labels	Replay	Shrink Perturb
1188							
1189	2015	0.700	0.700	0.700	0.704	0.700	0.700
1190	2016	0.705	0.703	0.707	0.708	0.701	0.707
1191	2017	0.702	0.706	0.709	0.709	0.703	0.707
1192	2018	0.707	0.709	0.711	0.711	0.700	0.706
1193	2019	0.707	0.708	0.708	0.710	0.701	0.704
1194	2020	0.703	0.709	0.709	0.710	0.704	0.705
1195	2021	0.708	0.707	0.711	0.711	0.708	0.707
1196	2022	0.705	0.709	0.710	0.709	0.705	0.706
1197							
1198							

Table 1: Area Under the Curve (AUC) performance on ProteinGym across different methods per year. Consistent with the Spearman correlation in Figure 5, Hare Tortoise achieves the strongest performance across all years, with Gradient Ascent and Shrink Perturb close behind.

## E UNIProt VALIDATION RESULTS

### E.1 UNIProt VALIDATION SET COMPOSITION

In this section, we provide additional information on the makeup of the validation set used Figure 4.

First, we show the taxonomic lineage breakdown of the validation set in Figure 16. We can see that the majority of sequences are Eukaryota, with some coverage of Bacteria and Archaea as well. The most common species in the validation set is *Homo sapiens* (human), followed by *Mus musculus* (house mouse) which makes sense given those are likely the most relevant species to drug discovery.

In Figure 17, we look at the pairwise sequence similarity between the different proteins in the validation set, and see that it is quite diverse, with the majority of sequence pairs being around 20-40% similar i.e. in the “protein twilight zone”. The “protein twilight zone” refers to the range of low sequence identity (typically 20-35%) where it becomes difficult to determine if two proteins are truly related based on their sequences alone.

### E.2 STRATIFICATION BY LINEAGE

In this section, we present the results on the UniProt validation set for each method stratified by the different lineages present in the dataset. We can see the results in Figure 18. Notably, the mean perplexity tends to follow the perplexity on Eukaryota and Archaea quite well, but Bacteria tends to have a much lower perplexity. Furthermore, as we go down the taxonomic tree, it does not seem to be the case that the model performs significantly better on the more common groups. There is a decent amount of spread in perplexity amongst the common groups, indicating that the models are not just memorizing the most common sequences.

## F FURTHER ABLATIONS

### F.1 WSD VS COSINE LEARNING RATE

In this section, we clarify the learning rate schedule used by our models. Our model is based off of AMPLIFY (Fournier et al., 2024), which used a cosine learning rate schedule in its training run. Unfortunately, because the cosine learning rate schedule has a fixed span it is unsuitable for continual training. Instead we use the warmup-stable-decay (WSD) schedule which has been used for continual pretraining (Li et al., 2025). In Figure 19, we can see that after decay, the two schedules perform about equivalently.

In our experiments, after each decay period, we reset to the checkpoint right before the decay before moving to the next task. Thus, only 90k out of the 100k gradient steps on a task are used to contribute to the continual training, but it offers a good balance between needing to decay the learning rate and being able to restart the run.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

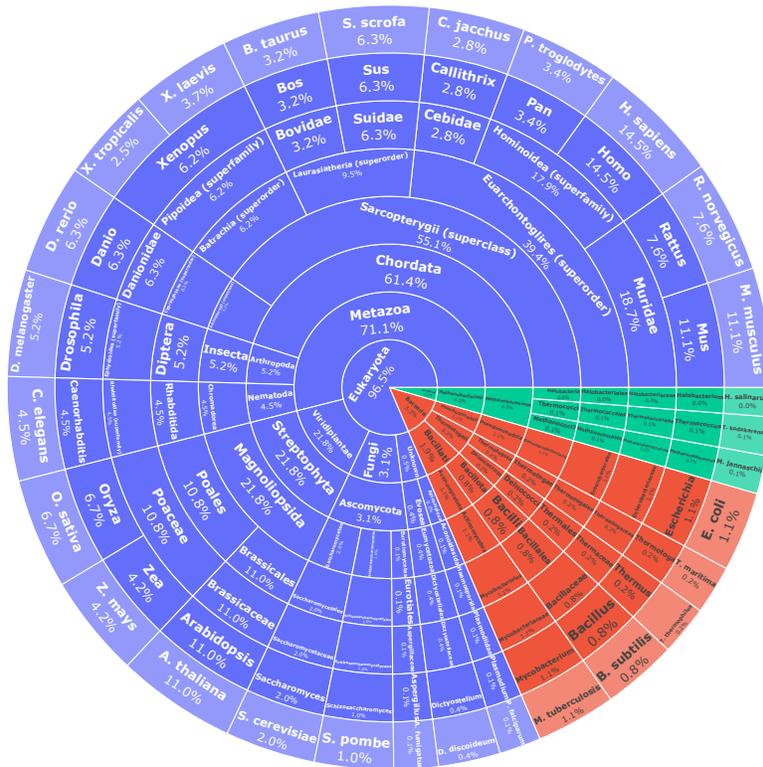


Figure 16: Proportion of different lineages in the UniProt validation set. Note in order to make Archaea properly visible, the area for each sector is according to the log of the number of sequences in that lineage.

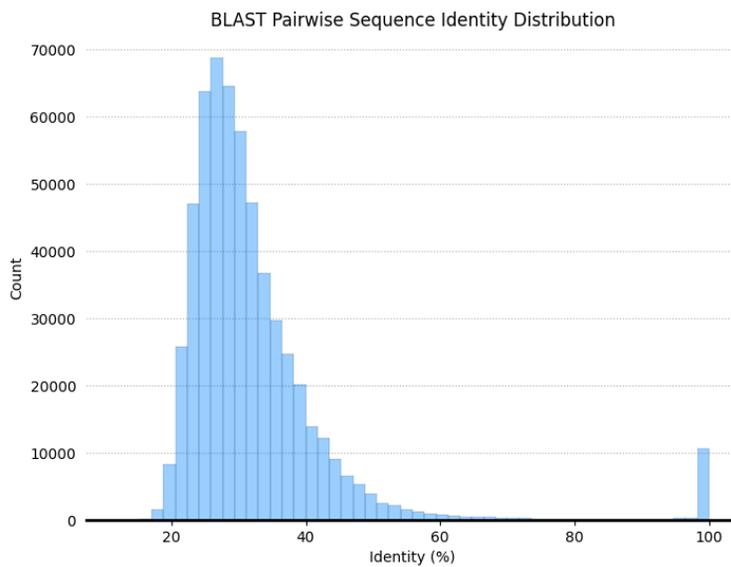


Figure 17: Sequence similarity between different proteins in the validation set.

1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

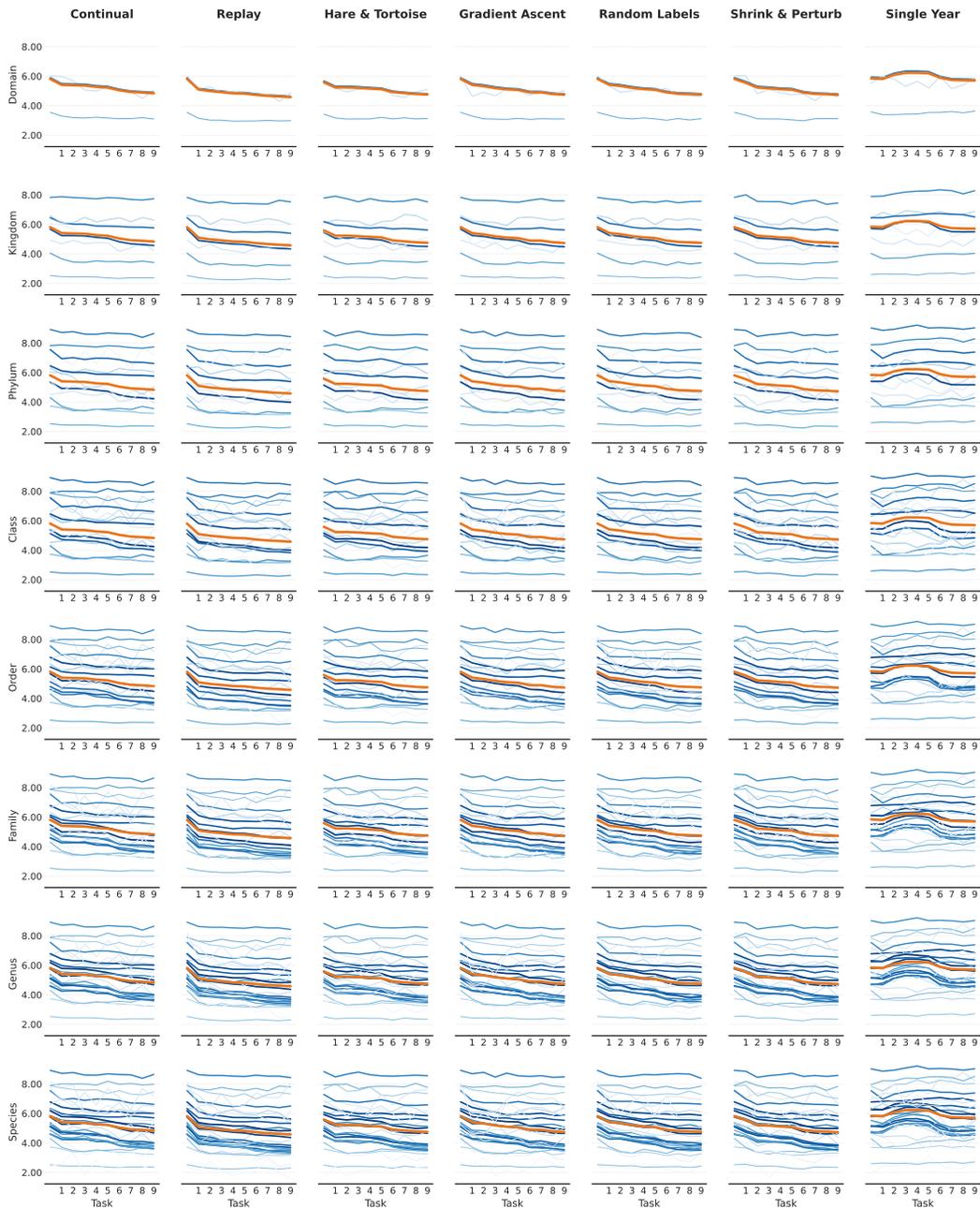


Figure 18: Perplexity on the UniProt validation set broken down by taxonomic lineage for each method. In each subplot, the mean perplexity is shown in orange, and the the lines for the more common groups are shaded to be darker and thicker.

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

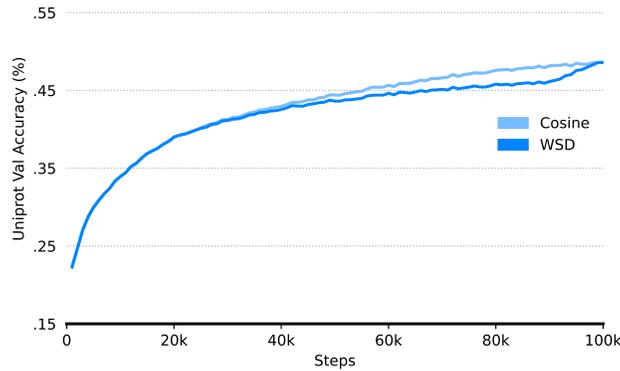


Figure 19: Both the cosine learning rate schedule and the warmup stable decay achieve approximately the same performance.

Method	UniProt Perplexity
2015 Data (910k steps)	4.535
Continual Sequence (910k steps)	4.568
Replay (910k steps)	<b>4.342</b>
AMPLIFY-1M (1 million steps)	4.359
Gradient Ascent (910k steps)	4.450
Hare and Tortoise (910k steps)	4.507
Random Labels (910k steps)	4.460
Shrink and Perturb (910k steps)	4.470

Table 2: Results of training for an equivalent number of on a single year (2015) compared to continual training across all years.

## F.2 LONGER TRAINING OF A SINGLE YEAR

In Table 2, we compare training for a longer period on a single year (2015) to continual training across all years. We find that training longer on a single year matches performance of continual training. This might be because the 2015 data is particularly representative of the overall UniProt distribution, as seen in Figure 3a. Regardless, all of the continual learning methods outperform the longer single year training, indicating that continual learning is beneficial beyond just training for longer.

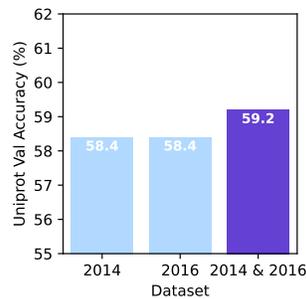


Figure 20: Data filtering experiments with a larger model (350M parameters).

Method	Hyperparameter	Distribution	Selected Value
Continual	None		
Replay	$\lambda_{replay}$	Uniform(0, 1.0)	0.357495045651384
Hare and Tortoise	$\lambda_{ht\_mom}$	Uniform(.5, 1.0)	.931247906596137
	$\lambda_{reset\_freq}$	LogInt(10, 10000)	559
Gradient Ascent	$\lambda_{asc}$	Uniform(0, 1.0)	0.0150798214665966
Random Labels	$\lambda_{rand}$	Uniform(0, 1.0)	0.00176366392582128
Shrink and Perturb	$\lambda_{shrink}$	Uniform(0, 0.9)	0.310430229773085
	$\lambda_{noise}$	Uniform(0, 1.0)	0.713412708958246

Table 3: The hyperparameter ranges and the selected hyperparameters for each method in our study.

Year	Release	Date	Number of Proteins
2015	2015_12	December 9, 2015	70,511,308
2016	2016_11	November 30, 2016	92,558,090
2017	2017_12	December 20, 2017	128,263,573
2018	2018_11	December 5, 2018	168,593,206
2019	2019_11	December 11, 2019	213,522,593
2020	2020_06	December 2, 2020	261,174,669
2021	2021_04	November 17, 2021	280,483,851
2022	2022_05	December 14, 2022	323,519,324
2023	2023_05	November 8, 2023	376,564,447
2024	2024_06	November 27, 2024	435,574,000

Table 4: The selected UniRef100 releases in our benchmark. The number of proteins listed are the numbers listed on the UniRef website, before we do any processing and deduplicating.

### F.3 DATA FILTERING EXPERIMENTS WITH LARGER MODELS

In Figure 20, we conduct a similar data filtering experiment as in Figure 3a, but with a larger model (350M parameters) in order to verify if the results hold at a larger scale or if they were an artifact of the smaller model potentially saturating performance. We select the intersection that had the best performance in the smaller model experiments (2014 intersected with 2016), and compare it to training on only 2014 data and only 2016 data. With the small model, we saw that training on the intersection outperformed training on either year alone. In Figure 20, we see the same trend, indicating that this is not an artifact of model size, and rather it is likely about the data quality itself.

## G HYPERPARAMETER SEARCH GRID

In Table 3, we describe the hyperparameter ranges and the selected value for each hyperparameter that was searched over in our study. Each hyperparameter was sampled independently, and we evaluated 16 trials for each method.

## H UNIREF STATISTICS

In Table 4, we list the specific releases we used to construct our benchmark.