

# HAIPR: A HIGH-THROUGHPUT AFFINITY PREDICTION FRAMEWORK FOR PROTEIN-PROTEIN INTERACTIONS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Accurate prediction of protein binding affinity is key for drug discovery and protein engineering, but commonly used evaluation protocols like Random Cross-Validation (RandomCV) can misrepresent true model generalization. We present HAIPR, a unified, open-source framework that streamlines the full machine learning pipeline for affinity prediction from training and optimization to inference, with curated benchmark datasets and robust, biologically meaningful evaluation protocols. By extending the BindingGYM benchmark and introducing realistic data splits, HAIPR reveals that RandomCV substantially overestimates model performance on out-of-distribution tasks. We systematically compare Support Vector Regression (SVR) using protein language model (pLM) embeddings to parameter-efficient fine-tuning (PEFT) of pLMs. SVR shows competitive results and increased stability in data-scarce scenarios, while PEFT excels as datasets grow larger and tasks become more complex. Analysis of model input setups shows that incorporating structural information does not always improve, and may sometimes hinder, performance for practical affinity prediction. Finally, we determine the lower limits of data required for reliable prediction, finding that even compact models can achieve performance close to the reproducibility limit of state-of-the-art assays, a practical ceiling for computational prediction. Code and pre-computed embeddings are publicly available.

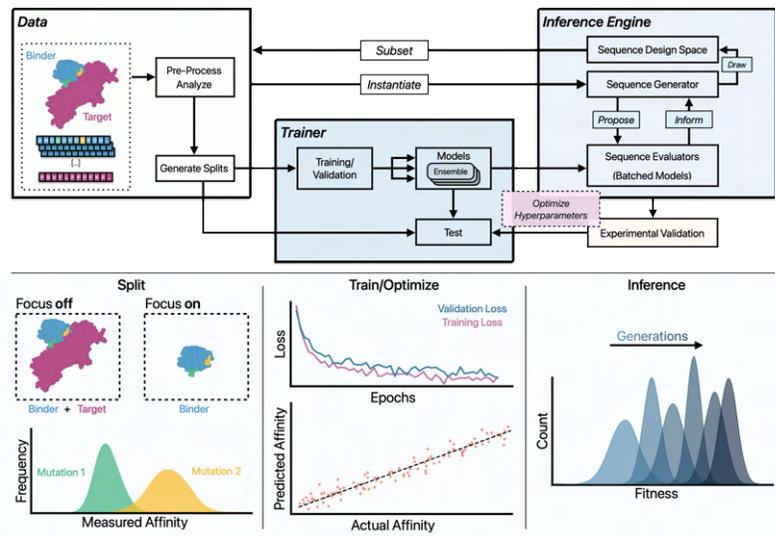


Figure 1: HAIPR Framework: We provide a unified framework for high-throughput affinity prediction that provides training, evaluation, and inference protocols as well as curated benchmark datasets.

## 1 INTRODUCTION

Protein-protein interactions (PPIs) are fundamental to cellular function, governing processes from signal transduction to immune responses Janin et al. (2008); Sprang (1997); Duc et al. (2015); Feinstein & Rowe (1965). Accurately predicting the effects of mutations on binding affinity in protein-protein complexes (PPCs) is a crucial step in drug development and protein engineering pipelines. Deep mutational scanning (DMS) has emerged as a powerful high-throughput screening (HTS) technique that enables systematic evaluation of thousands of mutations in parallel Adams et al. (2016), providing valuable data for machine learning approaches to predict binding affinity changes upon mutation Yang et al. (2019). Collections of DMS datasets were aggregated in benchmarks such as ProteinGym Notin et al. (2023) and BindingGYM Lu et al. (2024).

However, data availability remains a major bottleneck for developing foundation models for affinity prediction. Existing datasets such as SKEMPI2 Jankauskaitė et al. (2019) contain only few datapoints for any given complex, limiting the development of robust predictive models tailored to any specific complex. Furthermore, evaluation protocols commonly used in the literature, such as Random Cross-Validation (RandomCV), have been shown to overestimate model performance since train and test distributions are highly similar, leading to overly optimistic assessments of generalization capability Tossou et al. (2024).

Given these unknowns and challenges in the field, there is a clear need for a comprehensive framework that enables researchers to quickly evaluate various experimental setups, model architectures, and evaluation protocols in a standardized manner. Such a framework would facilitate unbiased comparisons across different algorithms and accelerate progress in the field.

Here, we address these challenges by introducing a comprehensive framework for high-throughput affinity prediction (HAIPR) that provides a unified evaluation protocol and interface to benchmark datasets. Our contributions are as follows:

1. We propose the HAIPR framework, which provides the backbone for evaluating future models by unifying the evaluation protocol and offering a unified access point to benchmark datasets, as well as comprehensive functionality to evaluate model performance and perturb input data.
2. We show that current splits to estimate out-of-distribution performance are insufficient, as they either fail to capture the true generalization challenges faced in real-world applications or use only a fraction of the available data, rendering them unsuitable for large-scale screenings.
3. We provide alternative splits to measure out-of-distribution performance: Leave-one-Mutation-out (LoMo) and Out-of-Distribution (OOD) splits that better reflect real-world generalization scenarios and utilize all available data.
4. We compare classical machine learning approaches such as Support Vector Regression (SVR) to parameter-efficient fine-tuning of protein language models (pLMs), demonstrating the relative strengths and limitations of each approach.
5. We evaluate the lower sample size threshold needed in DMS assays to achieve robust prediction performance, providing guidance for experimental design and data collection strategies.
6. We demonstrate the inference capabilities of the HAIPR framework to efficiently screen for novel variants that improve binding affinity based on fine-tuned pLMs trained on DMS data.

Our results demonstrate that while RandomCV can lead to overestimated performance, our proposed LoMo and OOD splits provide more realistic assessments of the generalization capabilities to unseen mutations and affinity ranges. We find that PEFT methods are more prone to model collapse but offer advantages when more data is available and especially when evaluating performance on OOD splits. Our analysis of data size requirements provides practical guidance for experimental design, showing that even relatively small datasets and models can achieve performance exceeding the resolution limits of DMS measurements given sufficiently similar training and test distributions.

## 2 RELATED WORK

Previous work has predominantly focused on general affinity prediction across diverse protein complexes, which differs fundamentally from our approach of learning single-complex scoring func-

108 tions. However, both the approaches often rely on pre-trained protein language models (pLMs) to  
109 provide embeddings of the protein sequences such as Evolutionary Scale Modeling (ESM) Lin et al.  
110 (2022), structural models such as ProteinMPNN Dauparas et al. (2022) or pLMs with additional  
111 structural context such as ESM-3 Hayes et al. (2025). This section reviews the relevant literature,  
112 highlighting the distinction between these two problem settings.

## 114 2.1 GENERAL BINDING AFFINITY PREDICTION OF PROTEIN-PROTEIN COMPLEXES

116 Predicting Binding Affinity changes upon mutation has been an active field of research for more than  
117 a decade Moretti et al. (2013). Early benchmarks such as SKEMPI Jankauskaitė et al. (2019) pro-  
118 vided datasets for evaluating mutation effects on binding affinities measured using low-throughput  
119 affinity assays, but were limited by small sample sizes for each complex. Liu et al. (2024a) extended  
120 this work by increasing the total sample size to 12157 by combining SKEMPI PDBbind Wang et al.  
121 (2005) and SabDAb Dunbar et al. (2014). More recent benchmarks such as ProteinGym Notin et al.  
122 (2023) and BindingGYM Lu et al. (2024) have addressed this limitation by providing a collection of  
123 preprocessed DMS datasets providing up to 92 thousand samples for a single complex and totaling  
124 up to half a million data points.

125 Many general predictors of binding affinity have been brought forward. Vangone & Bonvin (2015)  
126 demonstrated that interfacial contact networks can effectively predict binding affinity. Zhou et al.  
127 (2020) developed MuPIPR, an end-to-end deep learning framework that uses contextualized rep-  
128 resentations to estimate mutation effects on protein-protein interactions, achieving state-of-the-art  
129 performance on SKEMPI datasets. Fiorellini-Bernardis et al. (2024) proposed eGRAL, a graph neural  
130 network that combines ESM embeddings with structural information to predict binding affinity  
131 changes upon mutation. Jiao et al. (2025) demonstrated that pre-trained inverse folding models can  
132 effectively predict binding free energy changes ( $\Delta\Delta G$ ) for mutations in the SKEMPI dataset.

## 134 2.2 AFFINITY PREDICTION FOR SINGLE PROTEIN-PROTEIN COMPLEXES USING DMS DATA

136 Deep mutational scanning has emerged as a powerful experimental approach for high-throughput  
137 characterization of protein variants Moulana et al. (2022) Adams et al. (2016). These datasets are  
138 generated using high-throughput assays, generally relying on a combination of sorting and sequenc-  
139 ing as proposed by Adams et al. (2016). Although these measurements can contain systematic biases  
140 Trippe et al. (2022), their overall correlation with low-throughput assays can approach that of inter-  
141 assay correlation between low-throughput assays Kamat & Rafique (2017) Moulana et al. (2022).  
142 DMS typically tests all single point mutations of the scanned area, often the entire protein, leading  
143 to a large number of datapoints. DMS datasets enable a new approach to binding-affinity prediction  
144 by providing sufficient data to train complex-specific models. Jones & Thornton (1996) emphasized  
145 over two decades ago, that different types of protein-protein interactions may require tailored ap-  
146 proaches rather than one-size-fits-all models thus further supporting this approach. Kastiris et al.  
147 (2011) and Moal et al. (2011) also argued for complex-specific energy functions and highlighted the  
important trade-off between compute cost and prediction accuracy.

148 Lee et al. (2018) showed that deep mutational scanning data can predict evolutionary success,  
149 demonstrating the value of large-scale experimental data for training predictive models. Riesselman  
150 et al. (2018) demonstrated that deep generative models like DeepSequence can predict mutation ef-  
151 fects. Machine learning-guided protein engineering has shown remarkable success in optimizing  
152 protein functions with limited experimental data. Hie & Yang (2022) reviewed adaptive machine  
153 learning approaches for protein engineering, emphasizing sequential optimization strategies for dis-  
154 covering optimized sequences across multiple rounds of training and experimental measurement.

155 For antibody optimization, Bachas et al. (2022) developed deep learning approaches to predict both  
156 binding affinity and developability, enabling co-optimization of therapeutic antibodies. Shan et al.  
157 (2022) used geometric deep learning to optimize antibodies against SARS-CoV-2 variants, showcas-  
158 ing the potential for rapid in silico optimization. Gainza et al. (2023) used geometric deep learning  
159 frameworks to design novel protein binders, opening possibilities for designing binders for any  
160 target of interest. Bachas et al. (2022) demonstrated that deep contextual language models can  
161 quantitatively predict binding of antibody variants spanning three orders of magnitude in  $K_D$  range,  
revealing strong epistatic effects that highlight the need for intelligent screening approaches.

162 A critical challenge in high-throughput screening is ensuring model reliability when applied to novel  
163 variants. Dias & Kolaczowski (2017) highlighted the critical importance of data quality for train-  
164 ing accurate prediction models, suggesting that efforts should focus on curating high-quality, high-  
165 resolution datasets rather than simply developing more complex models.

166 Nevertheless, training models that can generalize to unseen mutations and affinity ranges still poses  
167 a challenge. This is highlighted by Tossou et al. (2024) who demonstrated the pitfalls of covariate  
168 shift in molecular interactions, and the resulting overestimation of model performance on random  
169 train/test splits.

170 While steps have been taken to address this challenge, oftentimes the resulting splitting mechanism  
171 discards the majority of the available data such as in the contig or modulo splits proposed by Notin  
172 et al. (2022; 2023). Fernandez-Diaz et al. (2024) introduced the AU-GOOD metric for evaluating  
173 model generalization, providing a framework for assessing model reliability on dissimilar proteins.  
174 Phillips et al. (2021) reconstructed binding affinity landscapes of five distinct SARS-CoV-2 Binding  
175 Partners (4 Antibodies and human ACE2). This work demonstrated how single mutations can carry  
176 much of the affinity variance for a given PPC.

177 While many predictors have been proposed for approximating sequence-function relationships using  
178 DMS data, only few have been experimentally tested in vitro. This might also be due to the lack of  
179 end-to-end pipelines for high-throughput screening. The only peer-reviewed work we are aware of  
180 that exercised high-throughput affinity screening based models trained on DMS data is Gelman et al.  
181 (2021) who trained an ensemble of convolutional-, graph-convolutional-, fully-connected neural  
182 networks and a linear regression model on one-hot encoded sequences to screen novel mutations.

### 184 3 METHODS

#### 185 3.1 DATA

186 We filtered the BindingGYM benchmark to datasets containing more than 3000 samples and ex-  
187 panded it with 5 datasets from Moulana et al. (2022), yielding a total of 21 PPCs. We extended the  
188 BindingGYM benchmark with additional datasets derived from combinatorial libraries. Combina-  
189 torial libraries are characterized by a high mean frequency of all mutations and a limited number  
190 of mutation sites. In contrast, most DMS datasets contain most mutations but with low frequency.  
191 These combinatorial datasets provide alternative evaluation protocols for out-of-distribution perfor-  
192 mance based on unseen mutations while not relying on single mutants. See Figure A.2.1 for details.

193 We compared two input regimes in line with Lu et al. (2024):

- 194 • **Focus-on:** only chains carrying variance are used. (mutated chains)
- 195 • **Focus-off:** The entire complex is used. (mutated and non-mutated chains alike.)

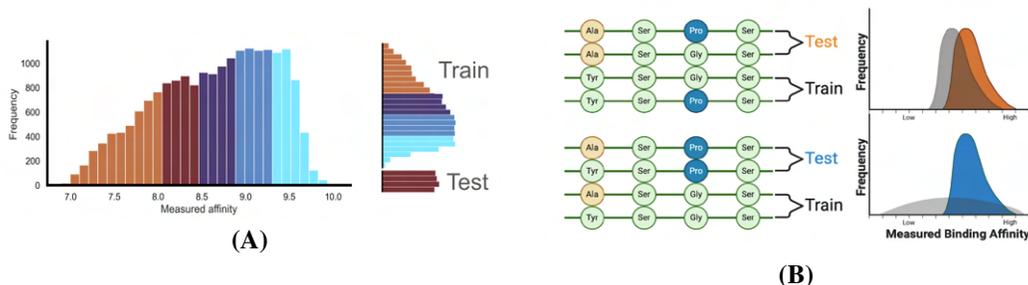
196 For datasets with a single mutated chain, we extracted that sequence and provided it to the (em-  
197 bedding) model. For datasets with more than one mutated chain or when we chose the "focus off"  
198 regime, we concatenated the sequences using a separator token. Sequences were tokenized with the  
199 model's native vocabulary. We obtained residue-level embeddings from the final layer and aggre-  
200 gated by mean pooling across the sequence length to obtain sequence embeddings unless specified  
201 otherwise.

#### 202 3.2 ALTERNATIVE SPLITS FOR OUT-OF-DISTRIBUTION ASSESSMENT

203 We proposed and evaluated two alternative splitting strategies that reflect real-world generalization  
204 challenges shown in Figure 2.

- 205 • **Out-of-Distribution (OOD) Split:** The target variable, in this case the affinity or affinity  
206 change, is divided into equally sized bins (based on sample count), with one bin held out  
207 for testing. This provides a more realistic assessment of model generalization to unseen  
208 affinity ranges while preserving a well balanced train to test ratio (Figure 2 Panel (A)).
- 209 • **Leave-One-Mutation-Out (LoMo) Split:** All sequences containing a specific mutation  
210 are excluded from training and used for testing. In a case of three mutable positions, each

216 having one alternative residue, this yields six possible splits, each containing half the total  
 217 available sequences. This approach directly measures how well the model predicts effects  
 218 of mutations that were absent during training. It is only applicable to combinatorial libraries  
 219 as shotgun screens would result in very few samples per split (Figure 2 Panel (B)).



222  
223  
224  
225  
226  
227  
228  
229  
230  
231 Figure 2: (A) Illustration of the Out-of-Distribution (OOD) split. Bins are computed to have equal  
 232 number of samples. (B) Illustration of the Leave-One-Mutation-Out (LoMo) split. All sequences  
 233 containing a particular mutation are removed from the training pool and used for testing. Example  
 234 splits contain all sequences with an Alanine at Position 1 (Orange; top panel), or all sequences that  
 235 have a Proline at Position 3 (Blue; bottom panel).

236  
237 These two approaches provide means to evaluate out-of-distribution performance while preserving  
 238 all available data. See Appendix Tables A.8.1 for the train test sample counts per benchmark and  
 239 split method.

### 240 241 3.3 MODELS

242 We evaluated two modeling approaches:

- 244 • **Support Vector Regression (SVR)**, using pLM embeddings as input features
- 245 • **Parameter-Efficient Fine-Tuning (PEFT)** of pLMs in combination with a simple regres-  
 246 sion head

247  
248 Within these approaches we probed various models from the ESM model family and ProteinMPNN  
 249 Dauparas et al. (2022). The model inputs can be either pre-computed embeddings (per residue or per  
 250 sequence) or the raw data, consisting of: (i) the protein sequence and (ii) optionally, the backbone  
 251 atom coordinates if structural information is used (i.e., `use_structure` is enabled). The specific  
 252 input type is determined by the configuration and model requirements.

253 The output for all our experiments is a scalar corresponding to the predicted fitness value (e.g., affini-  
 254 ty or affinity change). However, since output heads can be arbitrarily defined within the framework,  
 255 joint training and prediction of multiple target properties is also supported.

#### 257 258 3.3.1 SUPPORT VECTOR REGRESSION (SVR)

259 We fitted SVR models on pre-computed pLM features to probe the embedding space of the pLMs.  
 260 We used the scikit-learn implementation of SVR. If not stated otherwise, we instantiated the SVR  
 261 with the following hyperparameters:  $C=75$ ,  $\epsilon=0.1$ , `kernel=RBF`, and `gamma=scale`. We did not  
 262 constrain optimizer iterations but limited runtime to 48 CPU hours.

#### 263 264 3.3.2 PARAMETER-EFFICIENT FINE-TUNING (PEFT) OF PLMS

265 pLMs are trained on vast amounts of data which motivates their large parameter counts, but data  
 266 availability is often limited for downstream tasks. To reduce the number of trainable parameters  
 267 and adapt to the available dataset sizes, we employed parameter-efficient finetuning using the PEFT  
 268 library. We used weight decomposed low rank adaptation (DoRA) introduced by Liu et al. (2024b),  
 269 a matrix factorization approach that reduces the number of trainable parameters by factorizing the  
 weight matrix of the pLM into a low-rank approximation. In contrast to LoRa, DoRa factorizes

270 direction and magnitude separately, ensuring that the adapted weights are still on the unit sphere,  
271 thus closer resembling full fine-tuning Liu et al. (2024b).

272 For non-optimization runs we set rank to 2, alpha to 16, and dropout to 0.1. The MLP prediction  
273 head consisted of a single-layer MLP with hidden dimension size of 8, dropout of 0.5, and ReLU ac-  
274 tivation, mapping from the pLM embedding dimension to a single regression output. See Appendix  
275 Table A.8.3 for the resulting trainable parameters of the LoRa matrices and the MLP head as well  
276 as model abbreviations and sources.

## 277 278 279 4 RESULTS

### 280 281 4.1 THE HAIPR FRAMEWORK

282 We developed the HAIPR framework, which provides a unified backbone for evaluating affinity  
283 prediction models. HAIPR standardizes the evaluation protocol, offers a single access point to  
284 benchmark datasets, and includes comprehensive tools for model assessment and input data pertur-  
285 bation as well as inference. This design ensures consistency, reproducibility, and extensibility. The  
286 framework supports arbitrary models through a simple Predictor Interface. We provide new data  
287 splits while also supporting all splits from Notin et al. (2023), albeit arguing against the use of the  
288 latter. We support arbitrary sequence generators through a simple Generator Interface. We provide  
289 comprehensive customization of the framework through Hydra. We support optimization of most  
290 configurable parameters through Optuna enabling end-to-end optimization of all stages in unison.  
291 See Figure 1 for an overview of the framework.

### 292 293 4.2 LIMITATIONS OF RANDOM SPLITS AND CURRENT OUT-OF-DISTRIBUTION EVALUATION 294 PROTOCOLS

295 We trained SVR models on ESM embeddings of 21 large DSM datasets using RandomCV splits.  
296 For these datasets, we obtained mean Spearman correlation coefficients of 0.71 to 0.80 (Figure 3A).  
297 Larger models, such as ESM2-15B provided slightly better performance than smaller models. How-  
298 ever, for several datasets even the smallest ESM family model (ESM2-8M) achieved mean Spearman  
299 correlations on RandomCV that exceed the correlation between high-throughput and low-throughput  
300 assays and even sometimes between two distinct state-of-the-art low-throughput assays Kamat &  
301 Rafique (2017). This highlights the need for more realistic evaluation strategies. See Appendix A.3  
302 for a complete overview of the results.

303 One reason might be that individual mutations often account for a large portion of the variance in  
304 binding affinity. When using RandomCV, the model is exposed to these mutations during training,  
305 which could inflate performance estimates. However, in real-world deployment, the primary objec-  
306 tive is to accurately predict the effects of mutations that the model has not previously encountered,  
307 as shown previously in Tossou et al. (2024).

308 Alternative splitting strategies, such as Contig and Modulo splits Lu et al. (2024); Notin et al. (2023),  
309 are limited to single-mutation data and result in significant data loss, making them unsuitable for  
310 large-scale screenings (Figure 3 Panel (B)).

### 311 312 313 4.3 OOD AND LOMO SPLITS AS A SAMPLE EFFICIENT BIOLOGICALLY MOTIVATED 314 EVALUATION PROTOCOL

#### 315 316 4.3.1 LEAVE-ONE-MUTATION-OUT (LOMO)

317 To investigate the relationship between a models capability to predict an mutations absent in a  
318 combinatorial dataset and the mutations contribution to the variance in binding affinity, we trained  
319 SVR models using ESMC-300M embeddings on our proposed LoMo splits that withhold all samples  
320 containing a specific mutation from the training data. Figure 4 shows the results of the training  
321 using the splits that results in the highest and lowest mean affinity differences between train and test  
322 datasets for two different combinatorial assays introduced in Section 3.2. We found that mutations  
323 that shift the entire distribution of binding affinity pose a much greater challenge to the model than  
mutations that have only a small impact. Notably, the mutations evaluated are consistent across the

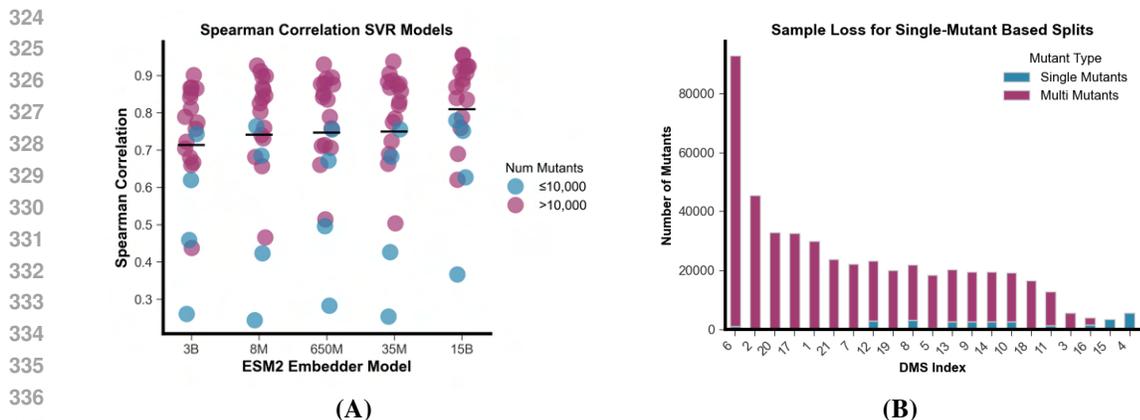


Figure 3: (A) Mean Spearman correlation over RandomCV of SVR models on ESM embeddings. (B) Number of samples split by single- (blue) and multi-mutants (red).

datasets, as they originate from the same library; however, some datasets contain fewer variants or lack certain mutations if those mutations prevented binding to the target.

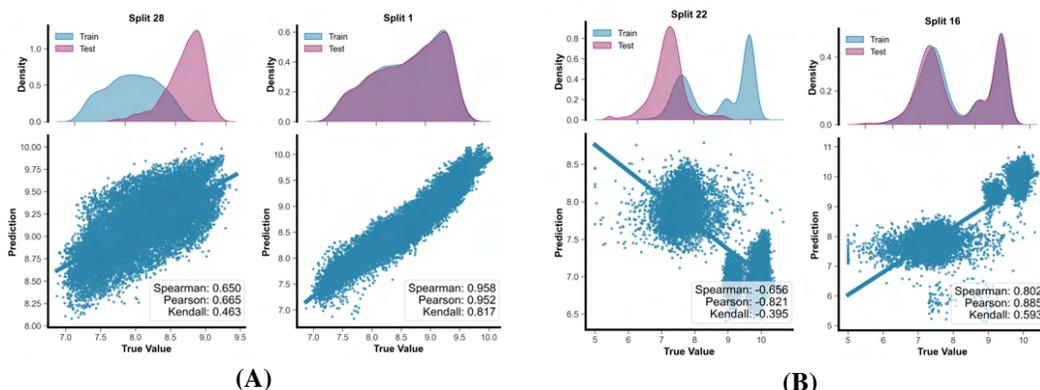


Figure 4: Using ESMC-300M Embeddings, we trained SVR models over LoMo splits for predicting the binding affinity between the SARS-CoV2 RBD and (A) Human ACE2 Receptor (DMS Index 17) as well as the (B) LY-CoV555 Antibody (DMS Index 19). Correlations and scatter plots are based on Test samples only while Density plots show the label distribution between train and test samples for the given split.

This further supports the argument that RandomCV is not a good proxy for high-throughput screening performance, where it is expected that the model is challenged by mutations that are not present in the training data. For a complete overview see Appendix A.5.

#### 4.3.2 OUT-OF-DISTRIBUTION (OOD)

We propose OOD splits as a general approach to evaluate out-of-distribution performance for high-throughput affinity prediction. This approach is less biologically motivated than LoMo splits where we specifically evaluate the model’s ability to predict unseen mutations but is not restricted to combinatorial libraries. To demonstrate the impact of our OOD splits on prediction performance, we trained SVR models on a range of ESM2 (8M, 150M, 650M, 15B) embeddings for all datasets containing fewer than 40,000 samples. Figure 5 shows the mean Spearman correlation for the CV and OOD splits. Comparison of model performance between RandomCV and OOD splits for SVR models trained on these ESM2 embeddings highlights the drop in predictive performance when moving from RandomCV splits, which overestimate generalization, to more realistic OOD splits that better reflect real-world scenarios.

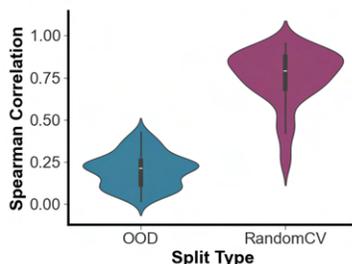


Figure 5: Distribution of Spearman correlation values for SVR models trained on ESM embeddings, comparing RandomCV and OOD splits for all benchmark datasets.

#### 4.4 EFFECTS OF STRUCTURAL CONTEXT ON MODEL PERFORMANCE

To evaluate the combined effect of using sequence only or sequence + structure information as well as providing the full Protein-Protein Complex in comparison to only the chains carrying any mutations, we trained SVR models on all combinations of focus and structure on four different ESM3 Embeddings, as ESM3 supports tokenization of structure information (Figure 6). While none of the reported metrics showed substantial differences the Focus-on without structure performed best overall, suggesting that additional data that does not carry variance is not helpful. Similarly, focus off did not improve performance for several tested sequence-only embeddings (ESM2: 8M, 35M, 650M; ESMC: 300M, 600M) across a subset of the BindingGYM benchmarks using OOD splits (Figure 18 in A.4.1. SVR models trained on OOD splits of ProteinMPNN embeddings (ca-only, v\_48\_020.pt) performed poorly (Table 27 in Appendix A.8.6). Thus providing additional invariant information does not substantially improve prediction performance.

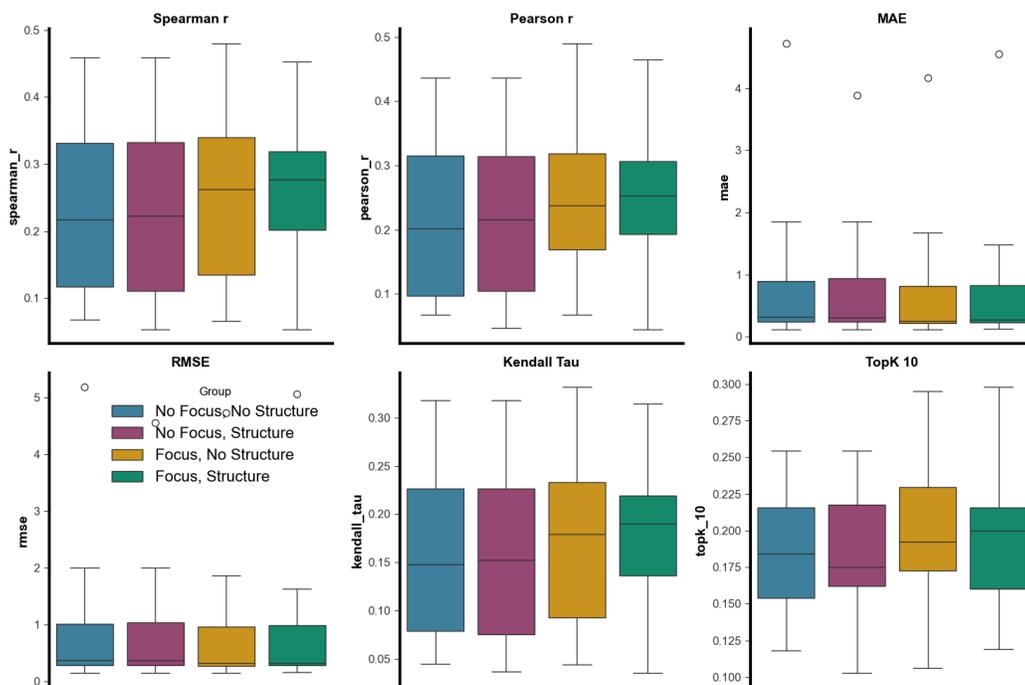


Figure 6: Mean Metrics over the full Benchmark for SVR models trained on ESM3 embeddings, comparing Focus-on and Focus-off as well as sequence only and sequence + structure.

## 4.5 SAMPLE SIZE REQUIREMENTS FOR RELIABLE PREDICTION

Next, we systematically investigated the minimum data size required to achieve reliable prediction performance. For benchmarks with more than 30,000 samples, we subsampled the training data at thresholds of 1,000, 2,500, 5,000, 10,000, 20,000 and 30,000 samples. Both SVR and PEFT models were evaluated using ESM2-8M, ESM2-15B, and ESMC-300M across OOD and CV splits. As expected, the effect of data size was more pronounced for OOD prediction, while even 1,000 samples were sufficient to exceed the data resolution for RandomCV prediction (Figure 7). While PEFT models outperformed SVR models, they were more challenging to train and in our experiments suffered more frequently from model collapse leading to missing datapoints. We hypothesize that further optimization of hyperparameters might mitigate this effect and we intend to explore this in future work. A graphic showing model collapse metrics can be found in Appendix A.6 Figure 31. Same Graphic with Root Mean Squared Error (RMSE) and Pearson  $r$  is in Appendix A.6 Figure 30.

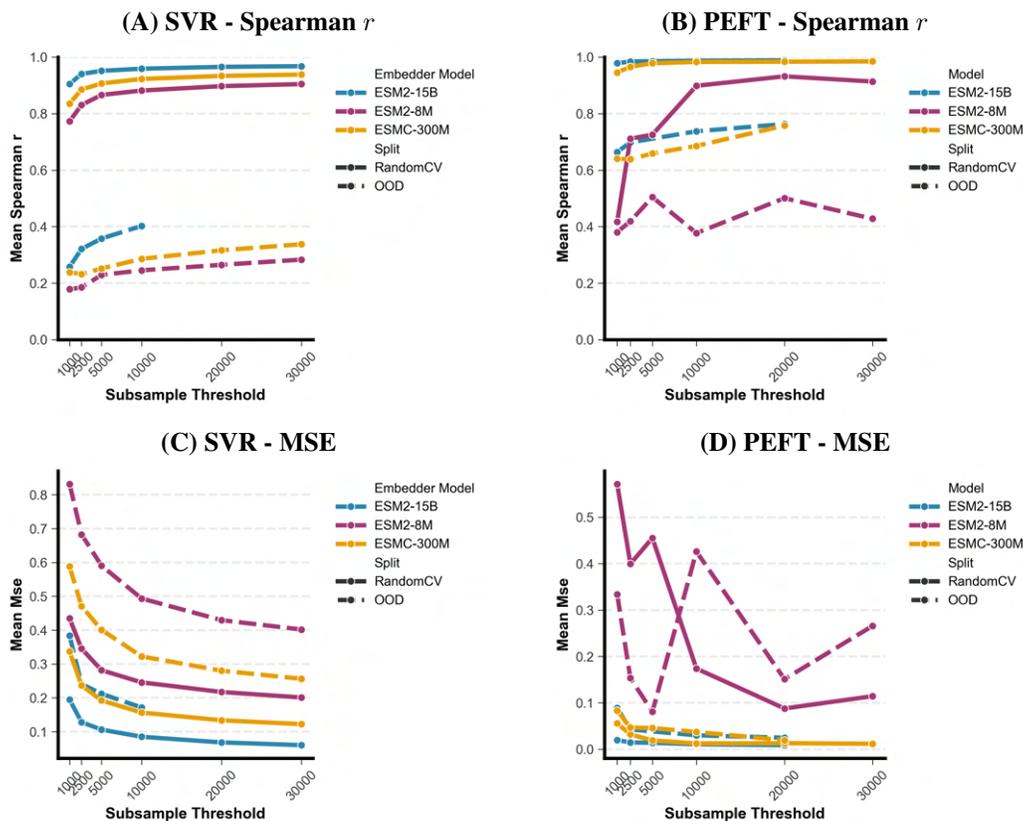
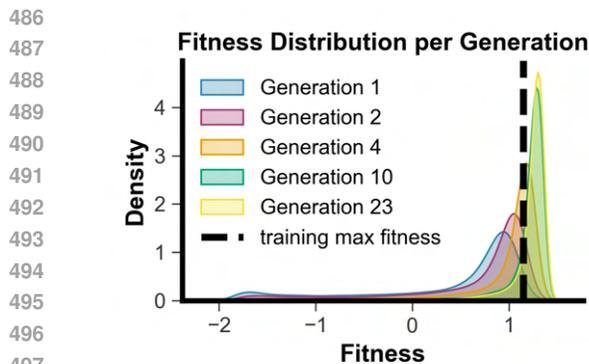


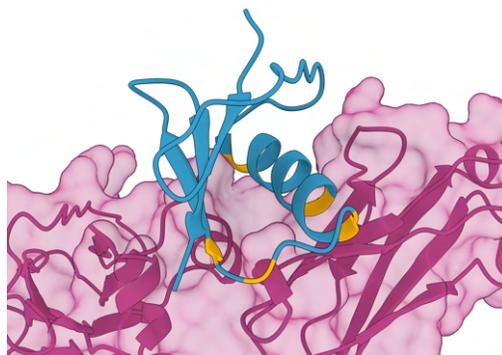
Figure 7: Sample size effects on predictive performance across model types and metrics. Panels show results for support vector regression (SVR, left) and PEFT models (right), evaluated using: (A,B) Spearman correlation and (C,D) Mean Squared Error (MSE) as a function of increasing training set size, across RandomCV and OOD splits. Missing data points indicate model collapse or cases exceeding a 48 hour training time limit.

## 4.6 HIGH-THROUGHPUT SCREENING AND DESIGN WITH HAIPR

We provide an implementation for sequence space exploration using a genetic algorithm based on Gad (2021). We used an ensemble of 5 ESMC-300M PEFT models trained on the GB1\_IgG-Fc\_fitness\_1FCC dataset on the OOD splits to score sequences generated by the genetic algorithm. Best generational sequences were folded using BOLTZ-2 Passaro et al. (2025) to ensure sequences still were predicted to fold (Figure 8 and Figure 9). For additional information see Appendix A.7 and Figure 35, 32 and 33 in the Appendix.



498  
499 Figure 8: Example of high-throughput in silico  
500 screening using HAIPR, showing predicted fit-  
501 ness scores for log-distributed generations.



502  
503 Figure 9: BOLTZ-2 prediction of the best gen-  
504 erational sequence, with mutations from wild-  
505 type highlighted in Orange.

## 506 5 DISCUSSION

507 The HAIPR framework offers streamlined, end-to-end access for high-throughput affinity predic-  
508 tion and inference using DMS assays. It emphasizes the importance of robust and realistic evalua-  
509 tion protocols, providing practical guidance for both experimental design and model selection. The  
510 framework is designed for easy extension to new models, facilitating rigorous and consistent eval-  
511 uation. Our findings highlight that evaluation splits encompassing the full affinity and mutational  
512 test distribution can significantly overestimate true out-of-distribution performance, underscoring  
513 the limitations of Random Cross-Validation in this context.

514 We have expanded the BindingGYM benchmark with five new combinatorial datasets, enabling  
515 sample-efficient assessment of out-of-distribution performance based on unseen mutations. These  
516 resources are now available to the community. While the introduced OOD and LoMo splits serve as  
517 a strong foundation for evaluating out-of-distribution generalization, we plan to further enhance the  
518 HAIPR framework with additional protocols for molecular OOD settings, such as those proposed  
519 by Fernandez-Diaz et al. (2024).

520 Our results show that SVR models struggle in challenging out-of-distribution settings but remain  
521 competitive with RandomCV and sufficient data. Unlike Loux et al., we found that incorporating  
522 structural information did not improve performance across broader benchmarks and tougher splits in  
523 Section 4.4. The data requirements for reliable prediction increase with evaluation split complexity,  
524 emphasizing the need for careful split design and providing practical guidance for designing DMS  
525 assays intended for training high-throughput predictors.

526 While hyperparameter optimization via Optuna was not a primary focus in this work, it is fully sup-  
527 ported by the framework. We provide examples for PEFT and training hyperparameters optimization  
528 in Appendix A.7.1 and A.8.7 demonstrating that even a limited number of trials can yield notable  
529 improvements in performance. Systematic exploration of hyperparameter optimization remains an  
530 avenue for future work. Initial inference experiments on the GB1\_IgG-Fc\_fitness\_1FCC dataset us-  
531 ing the OOD split produced promising results, which will be undergoing experimental validation in  
532 our laboratories.

533 Overall, HAIPR provides a unified interface for all presented experiments, enabling robust com-  
534 parison and fostering the development of future models for high-throughput Protein-Protein affinity  
535 prediction and design.

536  
537  
538  
539

## REFERENCES

- 540  
541  
542 Rhys M. Adams, Thierry Mora, Aleksandra M. Walczak, and Justin B. Kinney. Measuring the  
543 sequence-affinity landscape of antibodies with massively parallel titration curves. *eLife*, 5  
544 (DECEMBER2016), December 2016. ISSN 2050084X. doi: 10.7554/eLife.23156.
- 545 Sharrol Bachas, Goran Rakocevic, David Spencer, Anand V Sastry, Robel Haile, John M Sutton,  
546 George Kasun, Andrew Stachyra, Jahir M Gutierrez, Edriss Yassine, Borka Medjo, Vincent Blay,  
547 Christa Kohnert, Jennifer T Stanton, Alexander Brown, Nebojsa Tijanic, Cailen McCloskey, Re-  
548becca Viazzo, Rebecca Consbruck, Hayley Carter, Simon Levine, Shaheed Abdulhaqq, Jacob  
549Shaul, Abigail B Ventura, Randal S Olson, Engin Yapici, Joshua Meier, Sean McClain, Matthew  
550Weinstock, Gregory Hannum, Ariel Schwartz, Miles Gander, and Roberto Spreafico. Antibody  
551optimization enabled by artificial intelligence predictions of binding affinity and naturalness.  
5522022. doi: 10.1101/2022.08.16.504181.
- 553 J. Dauparas, I. Anishchenko, N. Bennett, H. Bai, R. J. Ragotte, L. F. Milles, B. I. M. Wicky,  
554 A. Courbet, R. J. de Haas, N. Bethel, P. J. Y. Leung, T. F. Huddy, S. Pellock, D. Tischer, F. Chan,  
555 B. Koepnick, H. Nguyen, A. Kang, B. Sankaran, A. K. Bera, N. P. King, D. Baker, R J De Haas,  
556 N. Bethel, P. J. Y. Leung, T. F. Huddy, S. Pellock, D. Tischer, F. Chan, B. Koepnick, H. Nguyen,  
557 A. Kang, B. Sankaran, A. K. Bera, N. P. King, and D. Baker. Robust deep learning based pro-  
558tein sequence design using ProteinMPNN. *bioRxiv*, pp. 2022.06.03.494563, June 2022. doi:  
55910.1101/2022.06.03.494563.
- 560 Raquel Dias and Bryan Kolaczowski. Improving the accuracy of high-throughput protein-protein  
561affinity prediction may require better training data. *BMC Bioinformatics*, 18(5):7–18, March  
5622017. ISSN 14712105. doi: 10.1186/S12859-017-1533-Z/FIGURES/4.
- 563 Nguyen Minh Duc, Hee Ryung Kim, and Ka Young Chung. Structural mechanism of G protein  
564activation by G protein-coupled receptor. *European Journal of Pharmacology*, 763:214–222,  
565September 2015. ISSN 0014-2999. doi: 10.1016/j.ejphar.2015.05.016.
- 566 James Dunbar, Konrad Krawczyk, Jinwoo Leem, Terry Baker, Angelika Fuchs, Guy Georges, Jiye  
567Shi, and Charlotte M. Deane. SAbDab: The structural antibody database. *Nucleic Acids Research*,  
56842(D1):D1140–D1146, January 2014. ISSN 0305-1048. doi: 10.1093/nar/gkt1043.
- 569 A. Feinstein and A. J. Rowe. Molecular Mechanism of Formation of an Antigen–Antibody  
570Complex. *Nature*, 205(4967):147–149, January 1965. ISSN 0028-0836, 1476-4687. doi:  
57110.1038/205147a0.
- 572 Raul Fernandez-Diaz, Hoang Thanh Lam, Vanessa López, and Denis C. Shields. A new frame-  
573work for evaluating model out-of-distribution generalisation for the biochemical domain. In *The  
574Thirteenth International Conference on Learning Representations*, October 2024.
- 575 Arturo Fiorellini-Bernardis, Sebastien Boyer, Christoph Brunken, Bakary Diallo, Karim Beguir,  
576Nicolas Lopez-Carranza, and Oliver Bent. Protein binding affinity prediction under multiple  
577substitutions applying eGNNs on Residue and Atomic graphs combined with Language model  
578information: eGRAL, May 2024.
- 579 Ahmed Fawzy Gad. PyGAD: An Intuitive Genetic Algorithm Python Library, June 2021.
- 580 Pablo Gainza, Sarah Wehrle, Alexandra Van Hall-Beauvais, Anthony Marchand, Andreas Scheck,  
581Zander Hartevelde, Stephen Buckley, Dongchun Ni, Shuguang Tan, Freyr Sverrisson, Casper  
582Goverde, Priscilla Turelli, Charlène Raclot, Alexandra Teslenko, Martin Pacesa, Stéphane Rosset,  
583Sandrine Georgeon, Jane Marsden, Aaron Petruzzella, Kefang Liu, Zepeng Xu, Yan Chai, Pu Han,  
584George F. Gao, Elisa Oricchio, Beat Fierz, Didier Trono, Henning Stahlberg, Michael Bronstein,  
585and Bruno E. Correia. De novo design of protein interactions with learned surface fingerprints.  
586*Nature*, 617(7959):176–184, May 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-05993-x.
- 587 Sam Gelman, Sarah A. Fahlberg, Pete Heinzelman, Philip A. Romero, and Anthony Gitter. Neural  
588networks to learn protein sequence–function relationships from deep mutational scanning data.  
589*Proceedings of the National Academy of Sciences*, 118(48):e2104878118, November 2021. ISSN  
5900027-8424, 1091-6490. doi: 10.1073/pnas.2104878118.

- 594 Thomas Hayes, Roshan Rao, Halil Akin, Nicholas J. Sofroniew, Deniz Oktay, Zeming Lin, Robert  
595 Verkuil, Vincent Q. Tran, Jonathan Deaton, Marius Wiggert, Rohil Badkundri, Irhum Shafkat,  
596 Jun Gong, Alexander Derry, Raul S. Molina, Neil Thomas, Yousuf A. Khan, Chetan Mishra,  
597 Carolyn Kim, Liam J. Bartie, Matthew Nemeth, Patrick D. Hsu, Tom Sercu, Salvatore Candido,  
598 and Alexander Rives. Simulating 500 million years of evolution with a language model. *Science*,  
599 387(6736):850–858, February 2025. doi: 10.1126/science.ads0018.
- 600 Brian L. Hie and Kevin K. Yang. Adaptive machine learning for protein engineering. *Current*  
601 *Opinion in Structural Biology*, 72:145–152, February 2022. ISSN 0959-440X. doi: 10.1016/J.  
602 SBI.2021.11.002.
- 603 J Janin, RP Bahadur, and P Chakrabarti Quarterly reviews of biophysics. Protein–protein interaction  
604 and quaternary structure. *cambridge.org*, 2008. doi: 10.1017/S0033583508004708.
- 605 Justina Jankauskaitė, Brian Jiménez-García, Justas Dapkūnas, Juan Fernández-Recio, and Iain H  
606 Moal. SKEMPI 2.0: An updated benchmark of changes in protein–protein binding energy, kinet-  
607 ics and thermodynamics upon mutation. *Bioinformatics*, 35(3):462–469, February 2019. ISSN  
608 1367-4803. doi: 10.1093/bioinformatics/bty635.
- 609 Xiaoran Jiao, Weian Mao, Wengong Jin, Peiyuan Yang, Hao Chen, and Chunhua Shen.  
610 BOLTSMANN-ALIGNED INVERSE FOLDING MODEL AS A PREDICTOR OF MUTA-  
611 TIONAL EFFECTS ON PROTEIN- PROTEIN INTERACTIONS. 2025.
- 612 Susan Jones and Janet M. Thornton. Principles of protein-protein interactions. *Proceedings of the*  
613 *National Academy of Sciences of the United States of America*, 93(1):13–20, January 1996. ISSN  
614 00278424. doi: 10.1073/PNAS.93.1.13.
- 615 Vishal Kamat and Ashique Rafique. Designing binding kinetic assay on the bio-layer interferometry  
616 (BLI) biosensor to characterize antibody-antigen interactions. *Analytical Biochemistry*, 536:16–  
617 31, November 2017. ISSN 0003-2697. doi: 10.1016/j.ab.2017.08.002.
- 618 Panagiotis L. Kastiris, Iain H. Moal, Howook Hwang, Zhiping Weng, Paul A. Bates, Alexan-  
619 dre M.J.J. Bonvin, and Joël Janin. A structure-based benchmark for protein-protein binding  
620 affinity. *Protein Science*, 20(3):482–491, March 2011. ISSN 09618368. doi: 10.1002/pro.580.
- 621 Juhye M. Lee, John Huddleston, Michael B. Doud, Kathryn A. Hooper, Nicholas C. Wu, Trevor Bed-  
622 ford, and Jesse D. Bloom. Deep mutational scanning of hemagglutinin helps predict evolutionary  
623 fates of human H3N2 influenza variants. *Proceedings of the National Academy of Sciences*, 115  
624 (35):E8276–E8285, August 2018. doi: 10.1073/pnas.1806133115.
- 625 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,  
626 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan Dos Santos Costa, Maryam Fazel-Zarandi,  
627 Tom Sercu, Salvatore Candido, and <sup>2</sup> Alexander Rives. Evolutionary-scale prediction of atomic  
628 level protein structure with a language model. *bioRxiv*, pp. 2022.07.20.500902, October 2022.  
629 doi: 10.1101/2022.07.20.500902.
- 630 Huaqing Liu, Peiyi Chen, Xiaochen Zhai, Ku-Geng Huo, Shuxian Zhou, Lanqing Han, and Guoxin  
631 Fan. PPB-Affinity: Protein-Protein Binding Affinity dataset for AI-based protein drug dis-  
632 covery. *Scientific Data*, 11(1):1316, December 2024a. ISSN 2052-4463. doi: 10.1038/  
633 s41597-024-03997-4.
- 634 Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-  
635 Ting Cheng, and Min-Hung Chen. DoRA: Weight-Decomposed Low-Rank Adaptation, July  
636 2024b.
- 637 Thomas Loux, Dianzhuo Wang, and Eugene I Shakhnovich. Does Structural Information Improve  
638 ESM3 for Protein Binding Affinity Prediction?
- 639 Wei Lu, Jixian Zhang, Ming Gu, and Shuangjia Zheng. BindingGYM: A Large-Scale Mutational  
640 Dataset Toward Deciphering Protein-Protein Interactions, December 2024.
- 641 Iain H. Moal, Rudi Agius, and Paul A. Bates. Protein-protein binding affinity prediction on a di-  
642 verse set of structures. 27(21):3002–3009, November 2011. doi: 10.1093/BIOINFORMATICS/  
643 BTR513.

- 648 Rocco Moretti, Sarel J. Fleishman, Rudi Agius, Mieczyslaw Torchala, Paul A. Bates, Panagio-  
649 tis L. Kastritis, João P.G.L.M. Rodrigues, Mikaël Trellet, Alexandre M.J.J. Bonvin, Meng Cui,  
650 Marianne Rooman, Dimitri Gillis, Yves Dehouck, Iain Moal, Miguel Romero-Durana, Laura  
651 Perez-Cano, Chiara Pallara, Brian Jimenez, Juan Fernandez-Recio, Samuel Flores, Michael  
652 Pacella, Krishna Praneeth Kilambi, Jeffrey J. Gray, Petr Popov, Sergei Grudin, Juan Esquivel-  
653 Rodríguez, Daisuke Kihara, Nan Zhao, Dmitry Korkin, Xiaolei Zhu, Omar N.A. Demerdash,  
654 Julie C. Mitchell, Eiji Kanamori, Yuko Tsuchiya, Haruki Nakamura, Hasup Lee, Hahnbeom Park,  
655 Chaok Seok, Jamica Sarmiento, Shide Liang, Shusuke Teraguchi, Daron M. Standley, Hiromitsu  
656 Shimoyama, Genki Terashi, Mayuko Takeda-Shitaka, Mitsuo Iwadate, Hideaki Umeyama, Dmitri  
657 Beglov, David R. Hall, Dima Kozakov, Sandor Vajda, Brian G. Pierce, Howook Hwang, Thom  
658 Vreven, Zhiping Weng, Yangyu Huang, Haotian Li, Xiufeng Yang, Xiaofeng Ji, Shiyong Liu,  
659 Yi Xiao, Martin Zacharias, Sanbo Qin, Huan Xiang Zhou, Sheng You Huang, Xiaoqin Zou,  
660 Sameer Velankar, Joël Janin, Shoshana J. Wodak, and David Baker. Community-wide evalua-  
661 tion of methods for predicting the effect of mutations on protein-protein interactions. *Proteins:  
662 Structure, Function and Bioinformatics*, 81(11):1980–1987, November 2013. ISSN 08873585.  
663 doi: 10.1002/prot.24356.
- 664 Alief Moulana, Thomas Dupic, Angela M. Phillips, Jeffrey Chang, Serafina Nieves, Anne A. Roffler,  
665 Allison J. Greaney, Tyler N. Starr, Jesse D. Bloom, and Michael M. Desai. Compensatory epistasis  
666 maintains ACE2 affinity in SARS-CoV-2 Omicron BA.1. *Nature Communications*, 13(1):7011,  
667 November 2022. ISSN 2041-1723. doi: 10.1038/s41467-022-34506-z.
- 668 Pascal Notin, Mafalda Dias, Jonathan Frazer, Javier Marchena-Hurtado, Aidan Gomez, Debora S.  
669 Marks, and Yarin Gal. Tranception: Protein fitness prediction with autoregressive transformers  
670 and inference-time retrieval, May 2022.
- 671 Pascal Notin, Aaron W Kollasch, Daniel Ritter, Lood van Niekerk, Steffanie Paul, Hansen Spinner,  
672 Nathan Rollins, Ada Shaw, Ruben Weitzman, Jonathan Frazer, Mafalda Dias, Dinko Franceschi,  
673 Rose Orenbuch, Yarin Gal, and Debora S Marks. ProteinGym: Large-Scale Benchmarks for  
674 Protein Fitness Prediction and Design. 2023.
- 675 Saro Passaro, Gabriele Corso, Jeremy Wohlwend, Mateo Reveiz, Stephan Thaler, Vignesh Ram  
676 Somnath, Noah Getz, Tally Portnoi, Julien Roy, Hannes Stark, David Kwabi-Addo, Dominique  
677 Beaini, Tommi Jaakkola, and Regina Barzilay. Boltz-2: Towards Accurate and Efficient Binding  
678 Affinity Prediction, June 2025.
- 680 Angela M Phillips, Katherine R Lawrence, Alief Moulana, Thomas Dupic, Jeffrey Chang, Milo S  
681 Johnson, Ivana Cvijovic, Thierry Mora, Aleksandra M Walczak, and Michael M Desai. Binding  
682 affinity landscapes constrain the evolution of broadly neutralizing anti-influenza antibodies. *eLife*,  
683 10:e71393, September 2021. ISSN 2050-084X. doi: 10.7554/eLife.71393.
- 684 Adam J. Riesselman, John B. Ingraham, and Debora S. Marks. Deep generative models of genetic  
685 variation capture the effects of mutations. *Nature Methods*, 15(10):816–822, October 2018. ISSN  
686 1548-7105. doi: 10.1038/s41592-018-0138-4.
- 687 Sisi Shan, Shitong Luo, Ziqing Yang, Junxian Hong, Yufeng Su, Fan Ding, Lili Fu, Chenyu Li,  
688 Peng Chen, Jianzhu Ma, Xuanling Shi, Qi Zhang, Bonnie Berger, Linqi Zhang, and Jian Peng.  
689 Deep learning guided optimization of human antibody against SARS-CoV-2 variants with broad  
690 neutralization. *Proceedings of the National Academy of Sciences*, 119(11):e2122954119, March  
691 2022. doi: 10.1073/pnas.2122954119.
- 692 Stephen R. Sprang. G PROTEIN MECHANISMS: Insights from Structural Analysis. *Annual Re-  
693 view of Biochemistry*, 66(Volume 66, 1997):639–678, July 1997. ISSN 0066-4154, 1545-4509.  
694 doi: 10.1146/annurev.biochem.66.1.639.
- 695 Prudencio Tossou, Cas Wognum, Michael Craig, Hadrien Mary, and Emmanuel Noutahi. Real-  
696 World Molecular Out-Of-Distribution: Specification and Investigation. *Journal of Chemical In-  
697 formation and Modeling*, 64(3):697–711, February 2024. ISSN 1549-9596. doi: 10.1021/acs.  
698 jcim.3c01774.
- 699 Brian L. Trippe, Buwei Huang, Erika A. DeBenedictis, Brian Coventry, Nicholas Bhattacharya,  
700 Kevin K. Yang, David Baker, and Lorin Crawford. Randomized gates eliminate bias in sort-seq  
701

702 assays. *Protein Science*, 31(9):e4401, September 2022. ISSN 1469-896X. doi: 10.1002/PRO.  
703 4401.  
704

705 Anna Vangone and Alexandre M.J.J. Bonvin. Contacts-based prediction of binding affinity in  
706 protein–protein complexes. *eLife*, 4(JULY2015), July 2015. ISSN 2050084X. doi: 10.7554/  
707 ELIFE.07454.

708 Renxiao Wang, Xueliang Fang, Yipin Lu, Chao-Yie Yang, and Shaomeng Wang. The PDBbind  
709 database: Methodologies and updates. *Journal of Medicinal Chemistry*, 48(12):4111–4119, June  
710 2005. ISSN 0022-2623. doi: 10.1021/jm048957q.  
711

712 Kevin K. Yang, Zachary Wu, and Frances H. Arnold. Machine-learning-guided directed evolution  
713 for protein engineering. *Nature Methods*, 16(8):687–694, July 2019. ISSN 15487105. doi:  
714 10.1038/s41592-019-0496-6.

715 Guangyu Zhou, Muhao Chen, Chelsea J.T. Ju, Zheng Wang, Jyun Yu Jiang, and Wei Wang. Mu-  
716 tation effect estimation on protein–protein interactions using deep contextualized representa-  
717 tion learning. *NAR Genomics and Bioinformatics*, 2(2), June 2020. ISSN 26319268. doi:  
718 10.1093/NARGAB/LQAA015.  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

A APPENDIX

A.1 IMPLEMENTATION DETAILS AND COMPUTE

All methods share a unified preprocessing and evaluation pipeline to ensure fairness across models and input regimes. We log experiments using the open source platform MLflow. We use the Optuna library for hyperparameter optimization. We intend to provide all generated embeddings to the community. All code used to conduct the experiments will be made available. If not stated otherwise we used the default parameters specified by the hydra configuration files.

A.2 FIGURES

A.2.1 COMBINATORIAL LIBRARIES

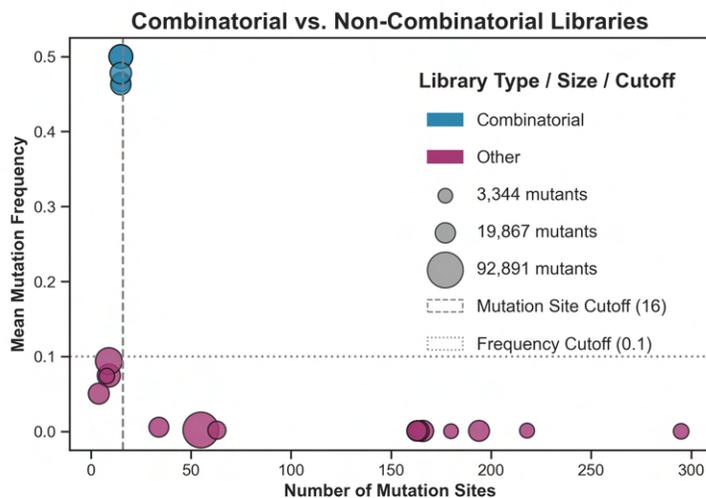


Figure 10: Combinatorial Libraries, are characterized by a high mean frequency of all mutations and a limited number of mutation sites.

A.3 ALL RESULTS FOR SVR-ESM2 MODELS ON RANDOM CV

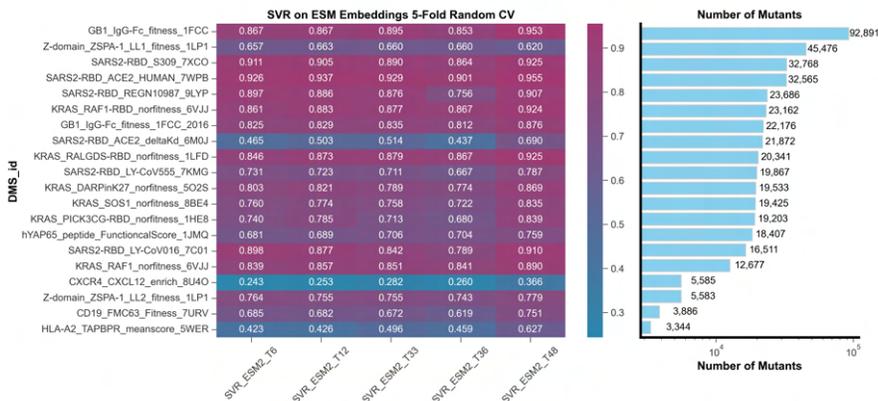


Figure 11: All Results for SVR-ESM2 Models on Random CV for spearman\_r

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

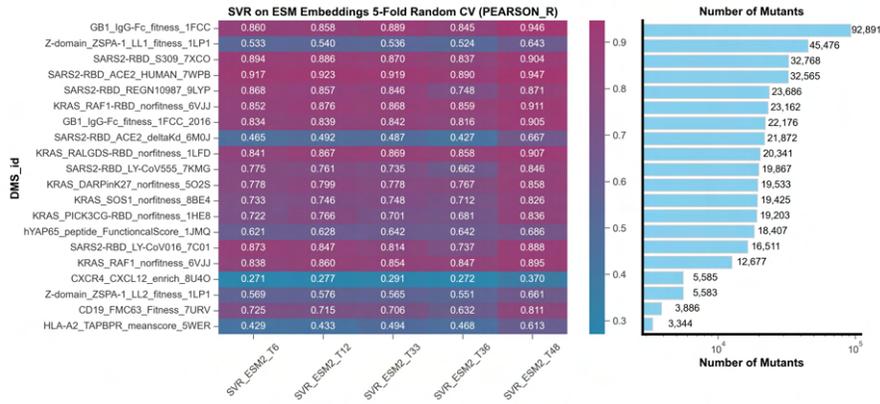


Figure 12: All Results for SVR-ESM2 Models on Random CV for pearson\_r

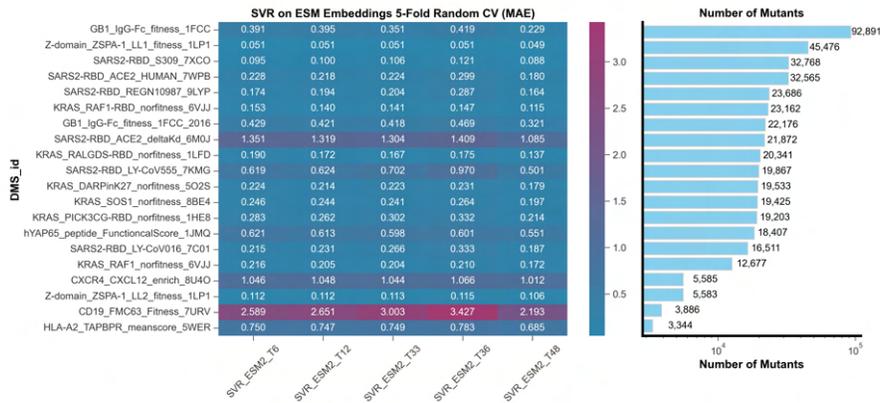


Figure 13: All Results for SVR-ESM2 Models on Random CV for mae

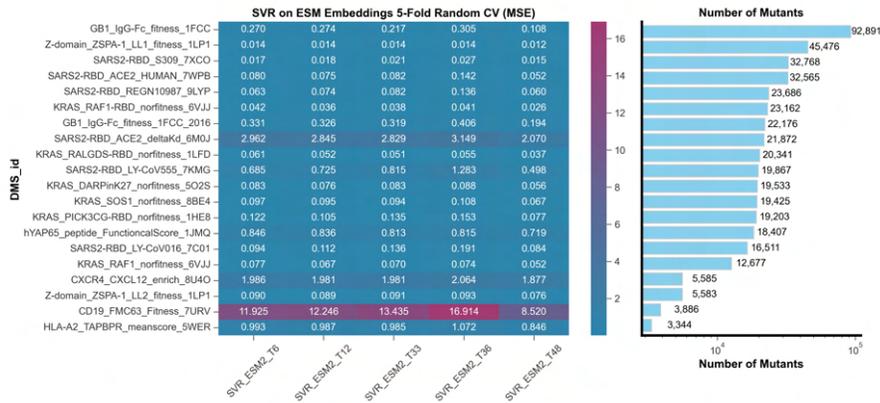


Figure 14: All Results for SVR-ESM2 Models on Random CV for mse

A.4 MLP vs PEFT

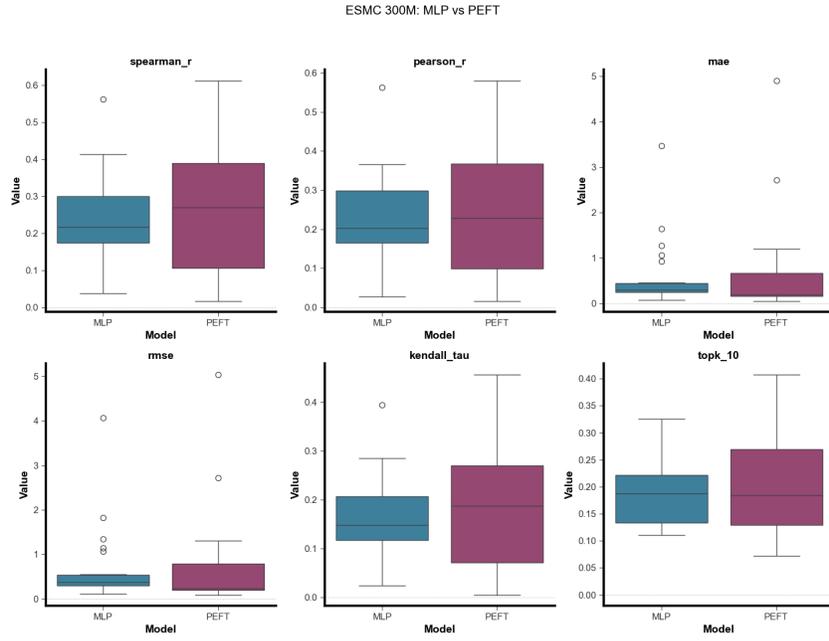


Figure 15: MLP vs PEFT boxplot difference for ESMC-300M

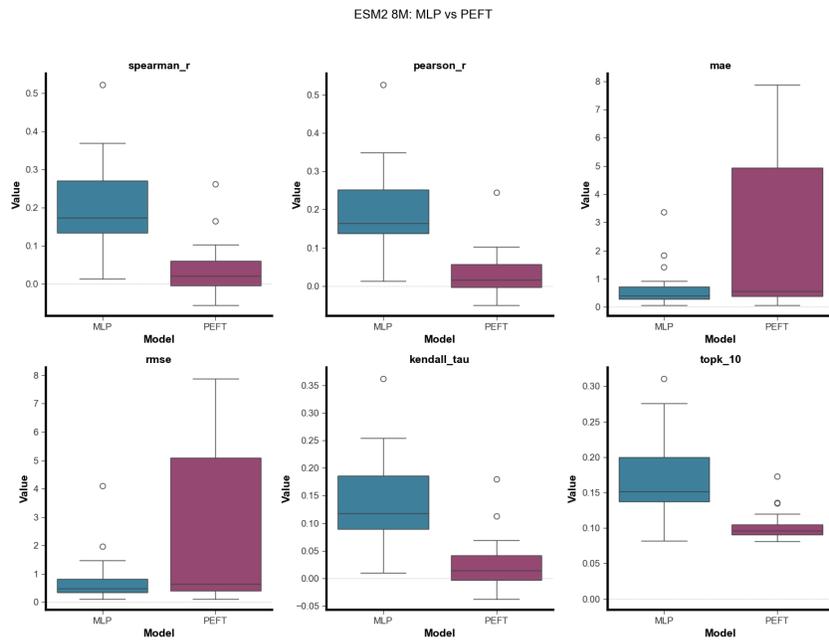


Figure 16: MLP vs PEFT boxplot difference for ESM2-8M

## A.4.1 FOCUS ON/OFF

To assess whether model would profit from providing full structural context while being sequence-only, we trained SVR models on ESM-family embeddings (ESM2: 8M, 35M, 650M; ESMC: 300M, 600M) across a subset of the BindingGYM benchmark using OOD splits. Figure 18 shows the distribution of Spearman correlation, comparing the two input regimes. Each violin plot represents the aggregated performance distribution for one regime over all datasets and models. The effect of providing full structural context was surprisingly small, although a small advantage of the Focus-on regimen was detectable. Future work will explore more direct approaches to investigate the impact on structural models.

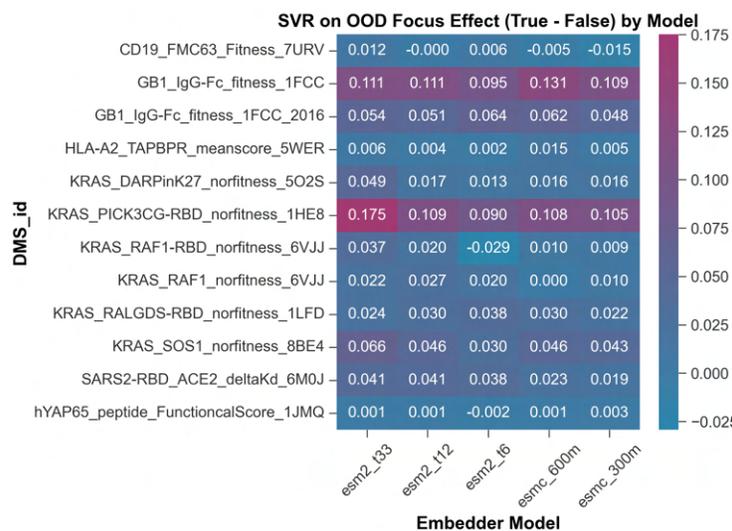


Figure 17: Results for BindingGYM subset comparing Support Vector Machines (SVR) on ESM embeddings with and without context (focus on/off)

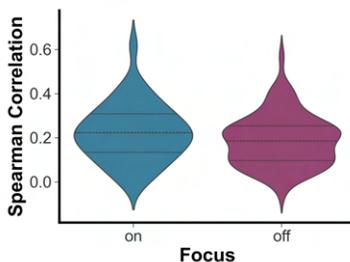


Figure 18: Distribution of Spearman correlation values for SVR models trained on ESM embeddings, comparing Focus-on and Focus-off.

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

## A.5 LOMO RESULTS

### A.5.1 MEAN SPEARMAN CORRELATION BY BENCHMARK AND MODEL FOR THE LOMO SPLIT

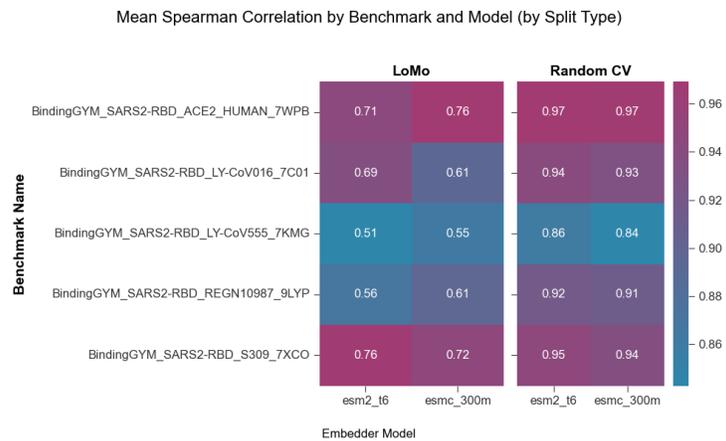


Figure 19: Mean Spearman correlation by benchmark and model for the LoMo split

### A.5.2 LOMO TRAIN-TEST LABEL DISTRIBUTION AND TEST-SET CORRELATION

The Splits numbers are based on mutation position in ascending order.

1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071  
 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079

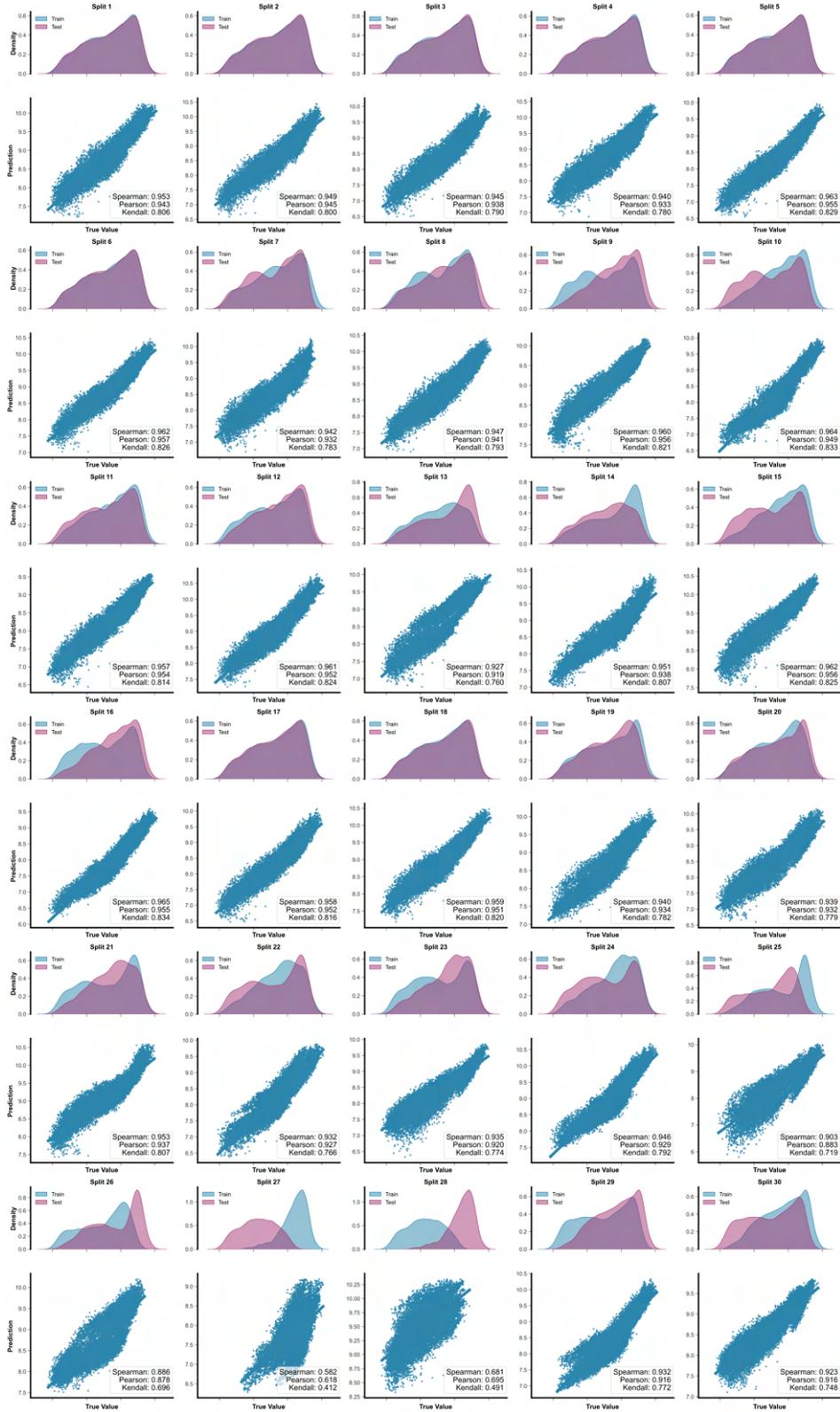


Figure 20: LoMo Train-Test label distribution and test-set correlation for ACE2 ESM2-8M

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

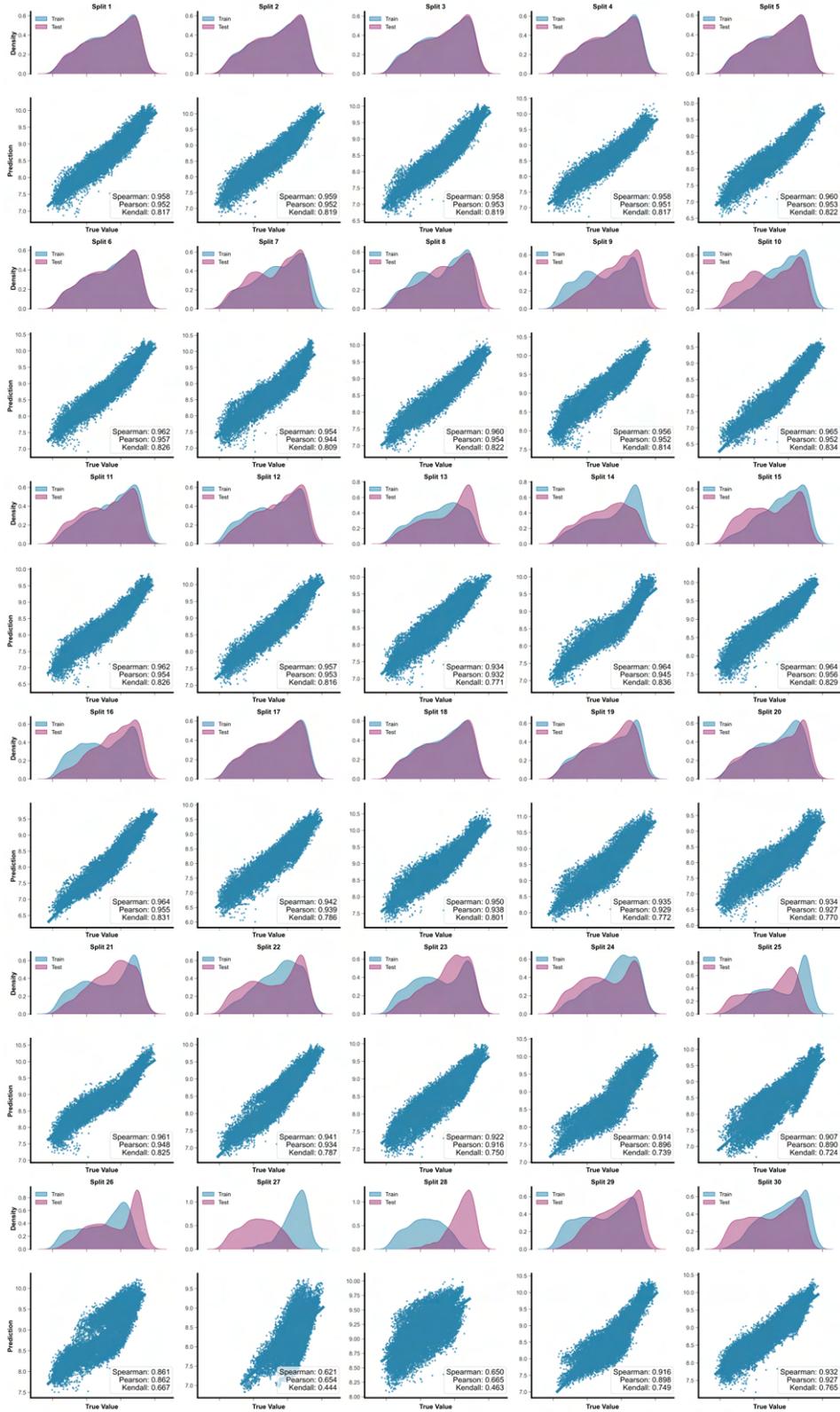


Figure 21: LoMo Train-Test label distribution and test-set correlation for ACE2 ESMC-300M

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

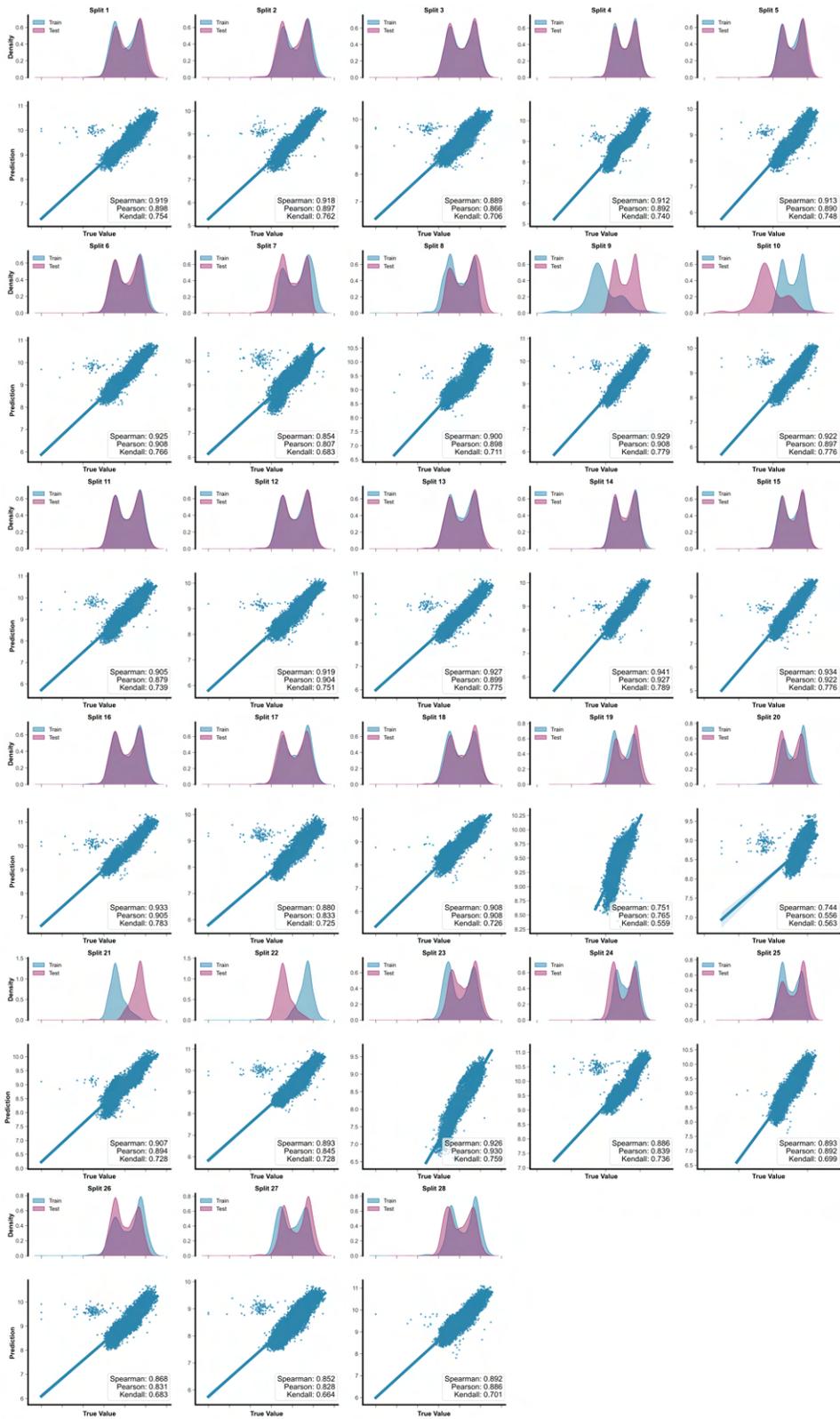


Figure 22: LoMo Train-Test label distribution and test-set correlation for LY-COV016 ESM2-8M

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

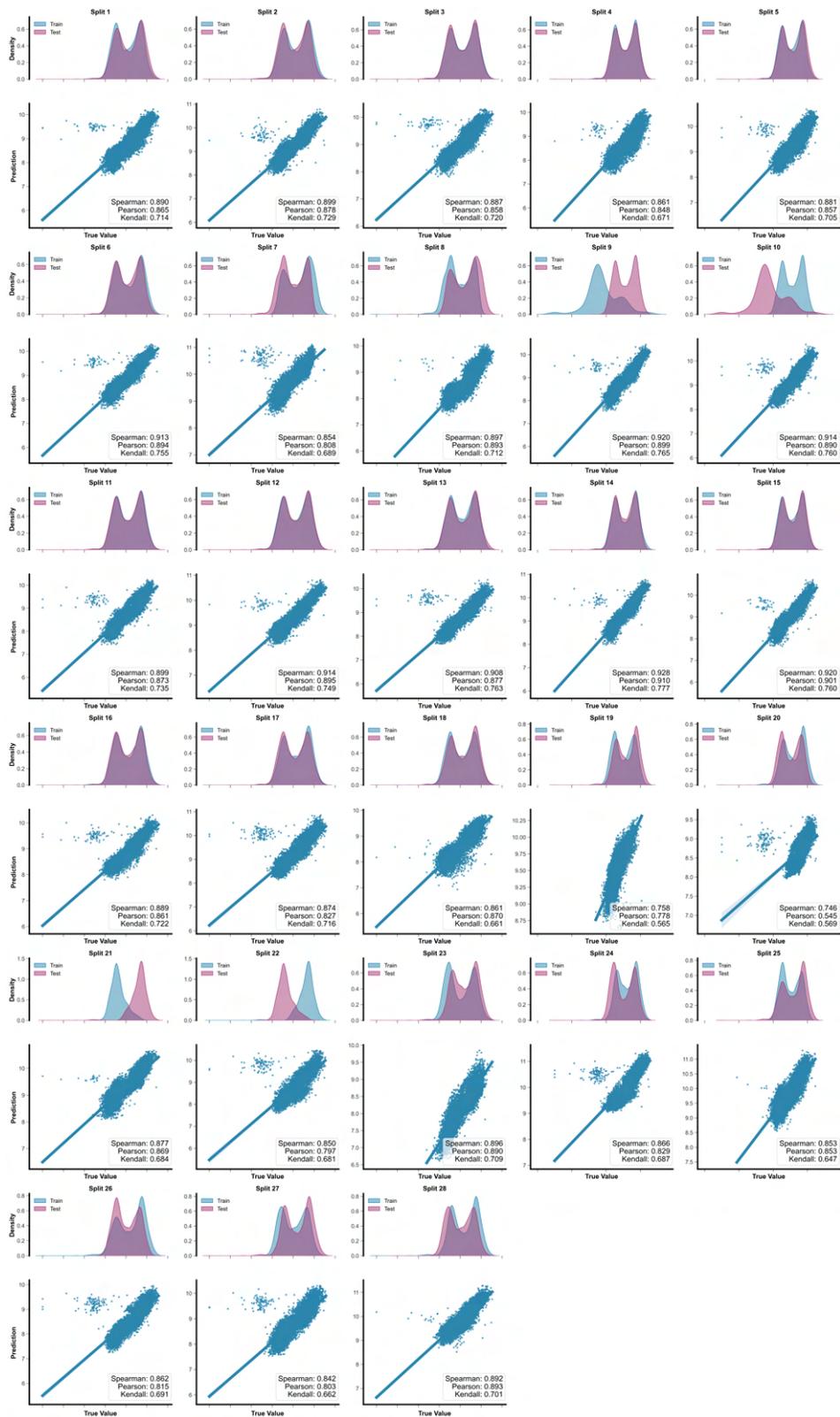


Figure 23: LoMo Train-Test label distribution and test-set correlation for LY-COV016 ESMC-300M

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

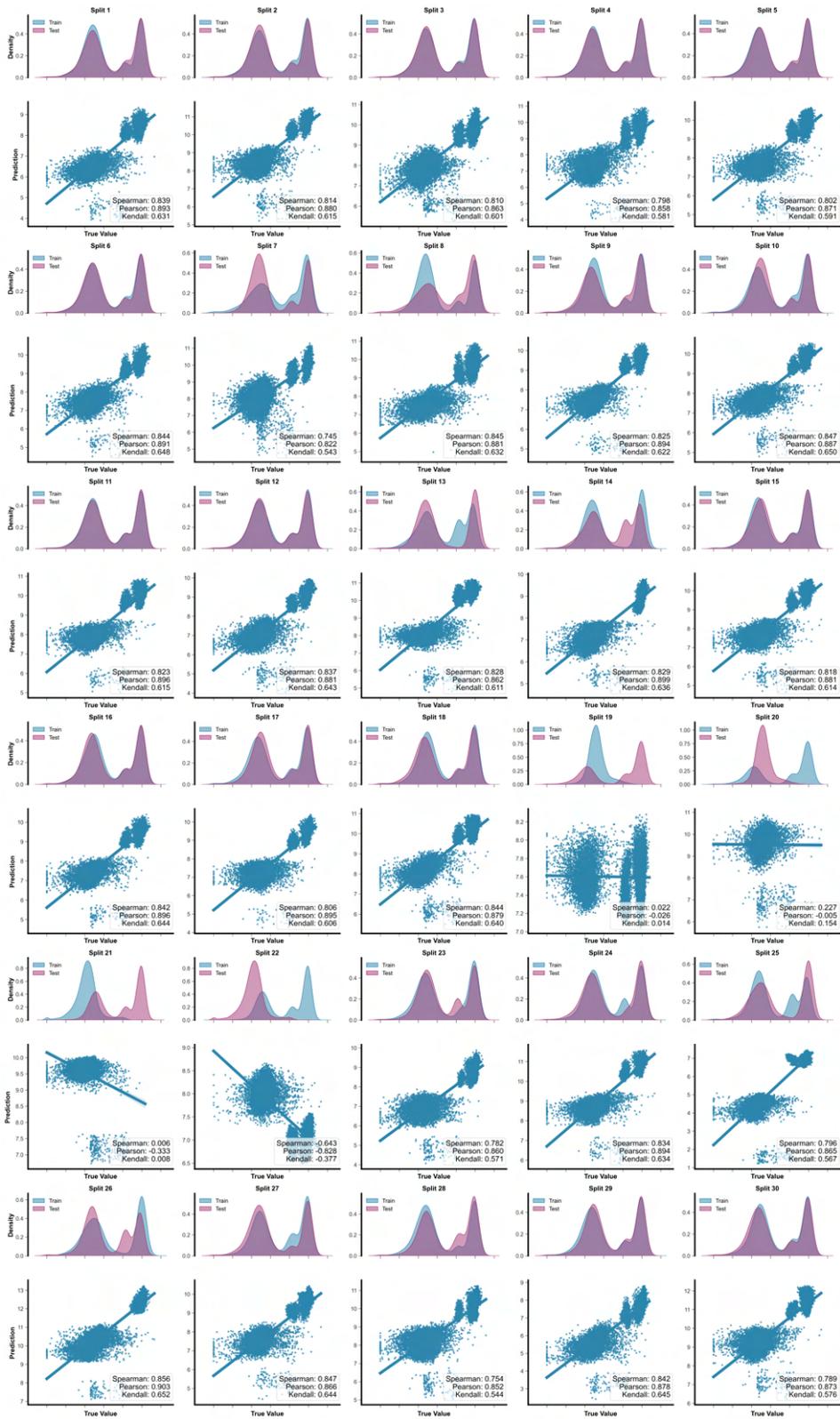


Figure 24: LoMo Train-Test label distribution and test-set correlation for LY-COV555 ESM2-8M

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

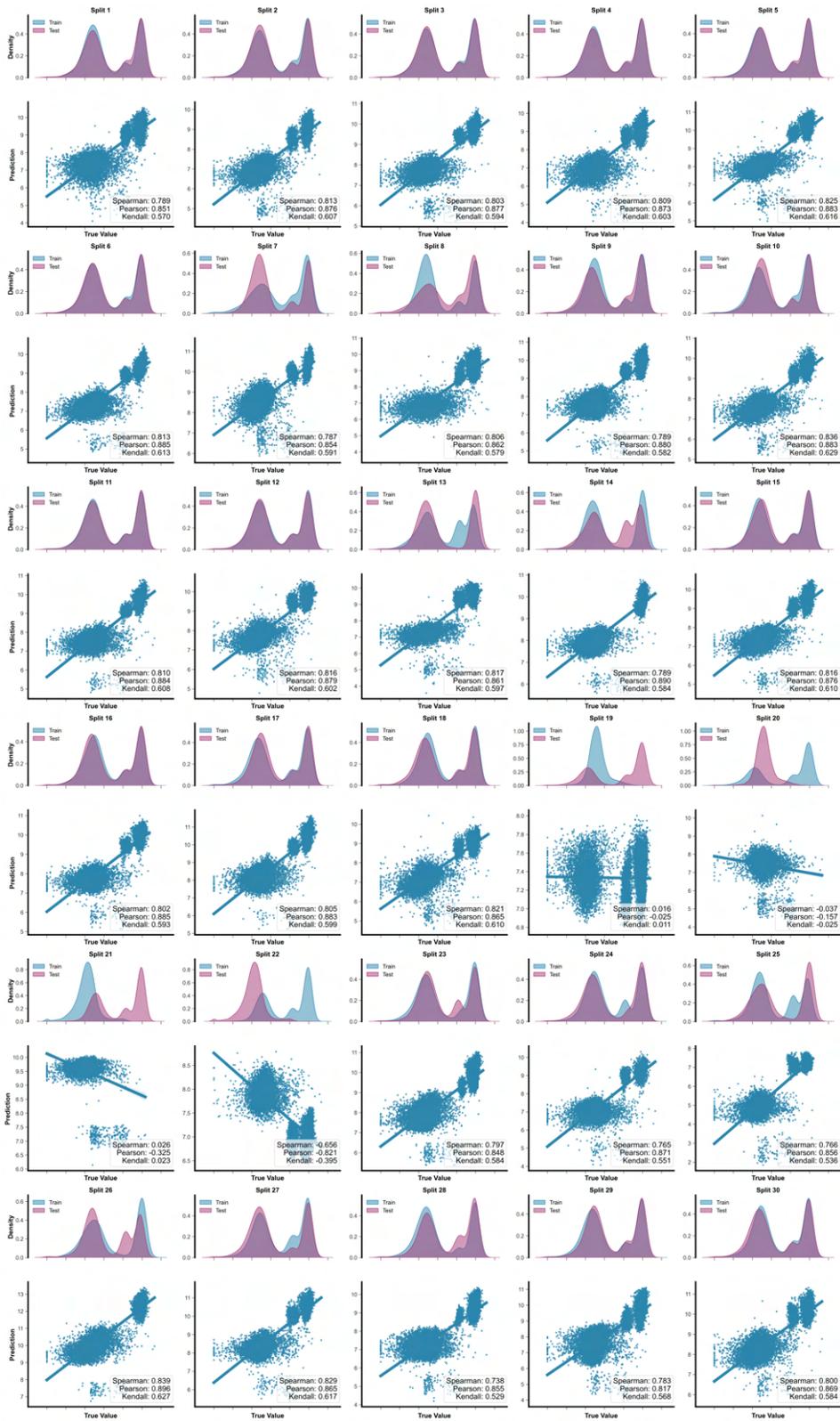


Figure 25: LoMo Train-Test label distribution and test-set correlation for LY-CoV555 ESMC-300M

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

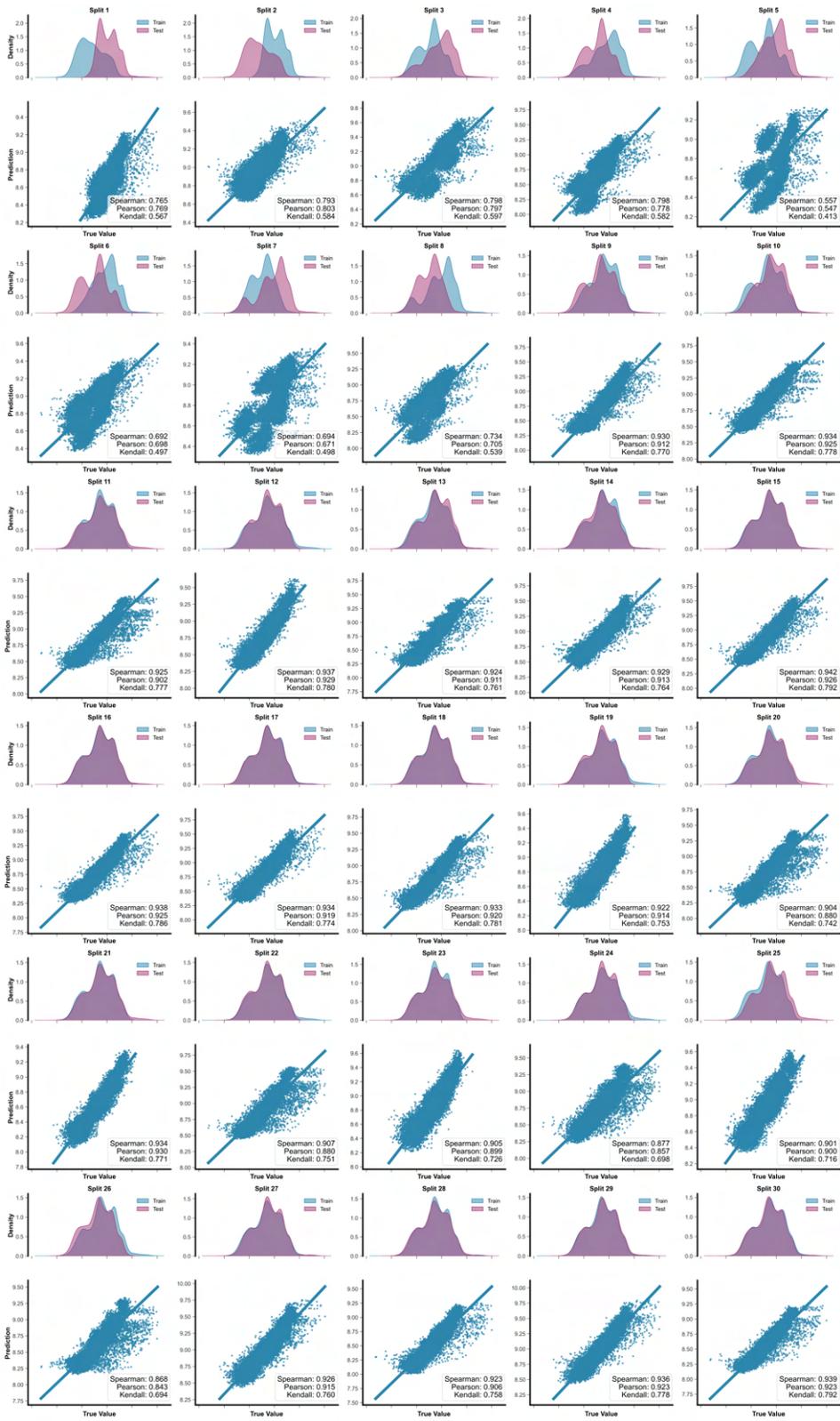


Figure 26: LoMo Train-Test label distribution and test-set correlation for S309 ESM2-8M

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

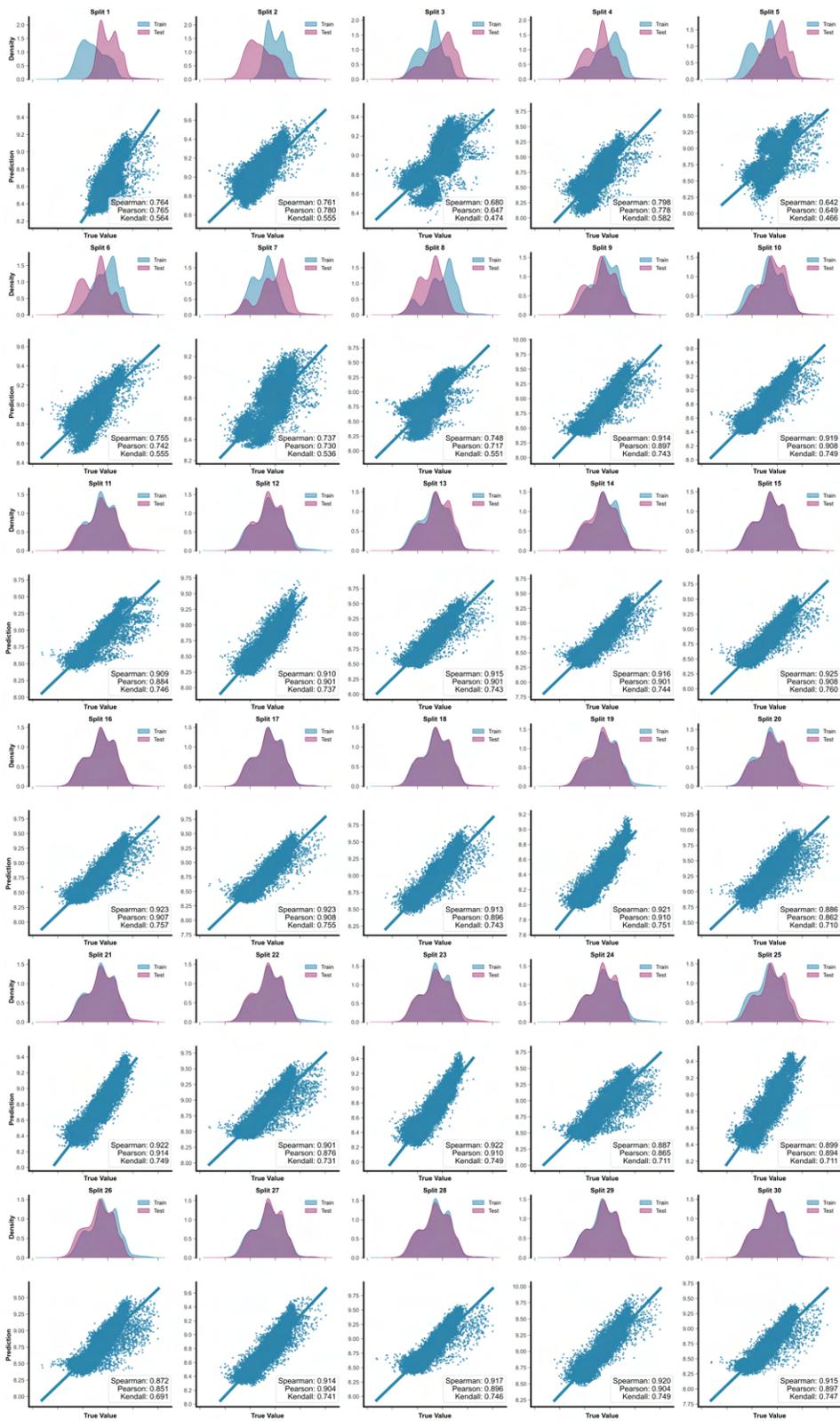


Figure 27: LoMo Train-Test label distribution and test-set correlation for S309 ESMC-300M

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

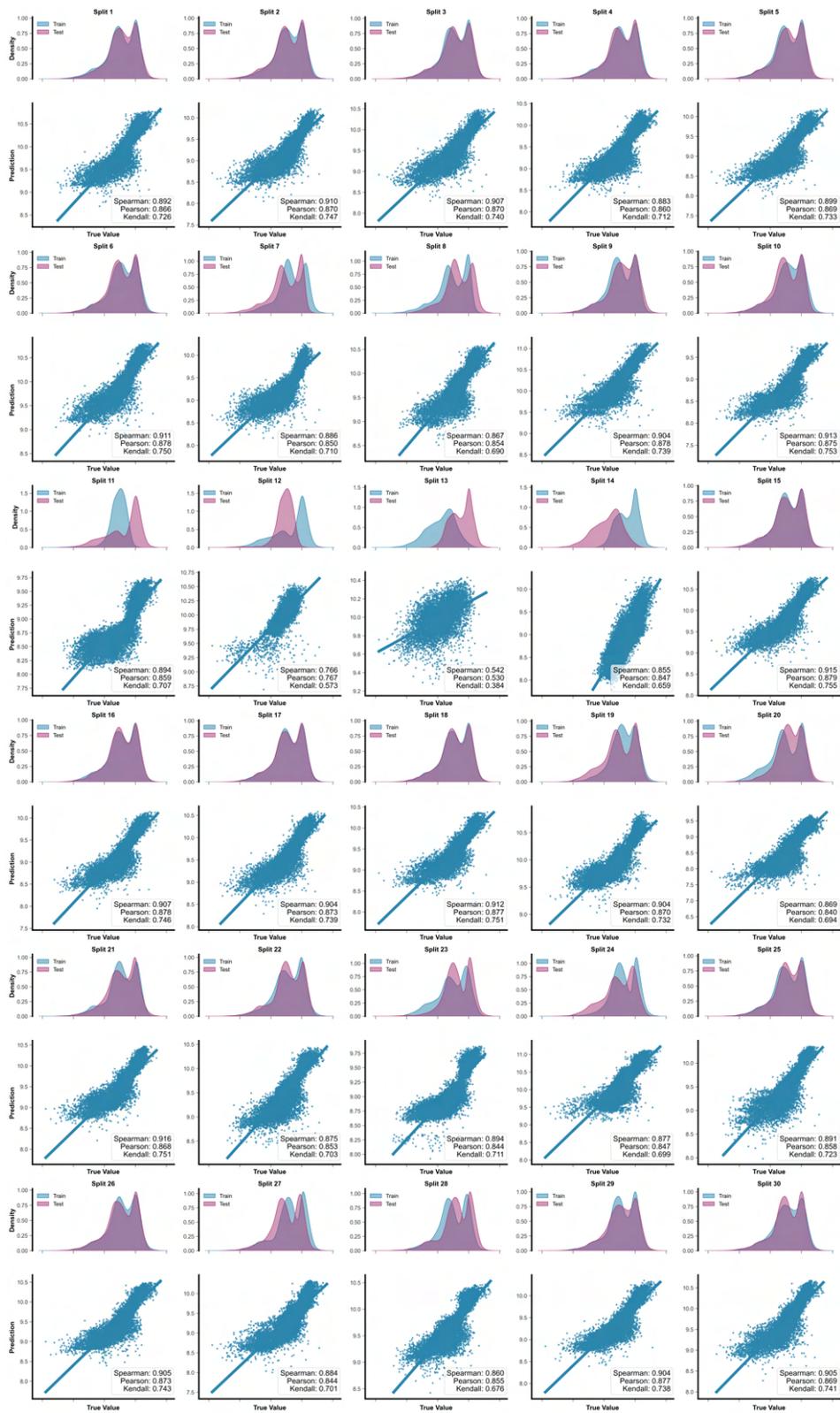


Figure 28: LoMo Train-Test label distribution and test-set correlation for REGN ESM2-8M

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

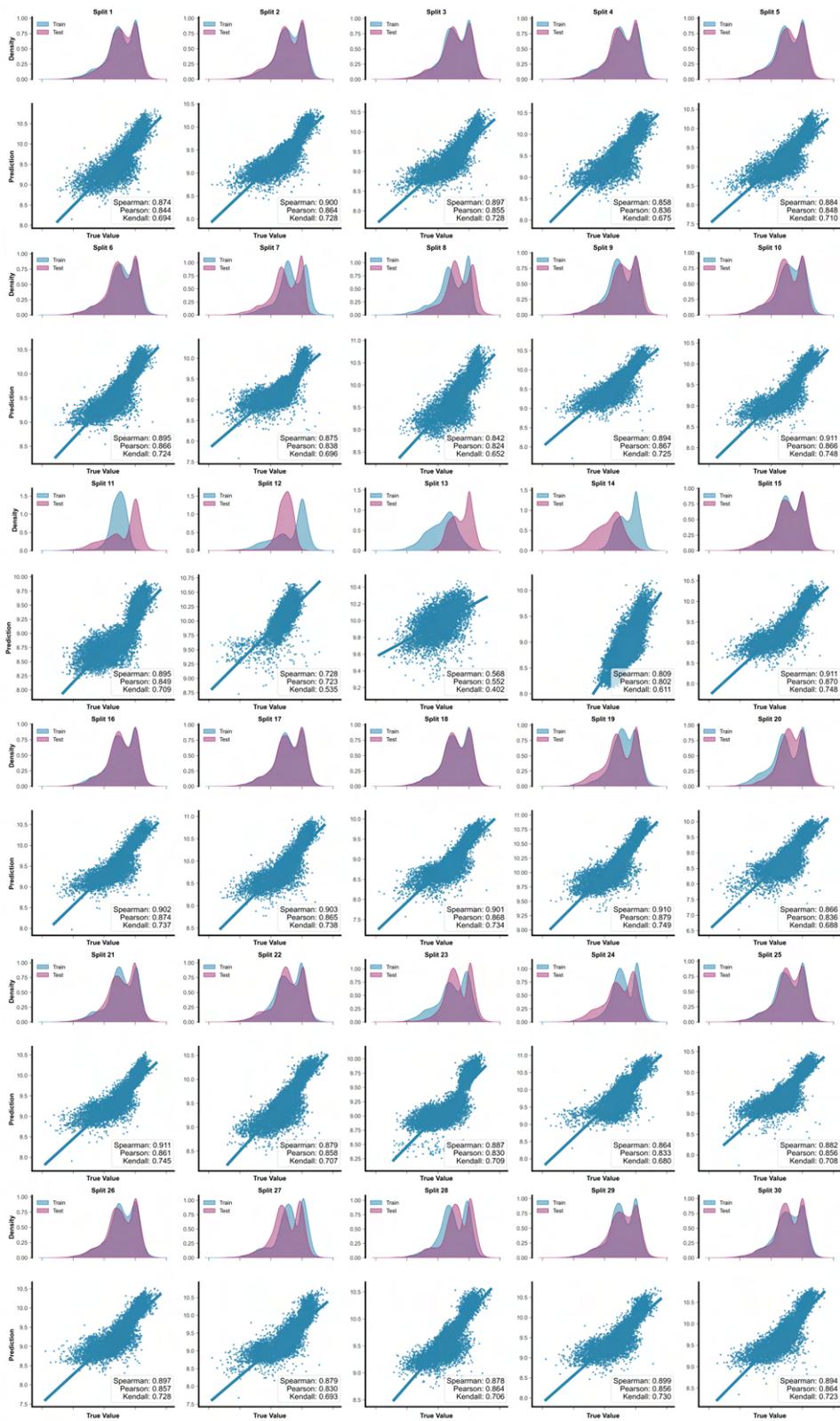


Figure 29: LoMo Train-Test label distribution and test-set correlation for REGN ESMC-300M

A.6 SAMPLE SIZE REQUIREMENTS FOR RELIABLE PREDICTION

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

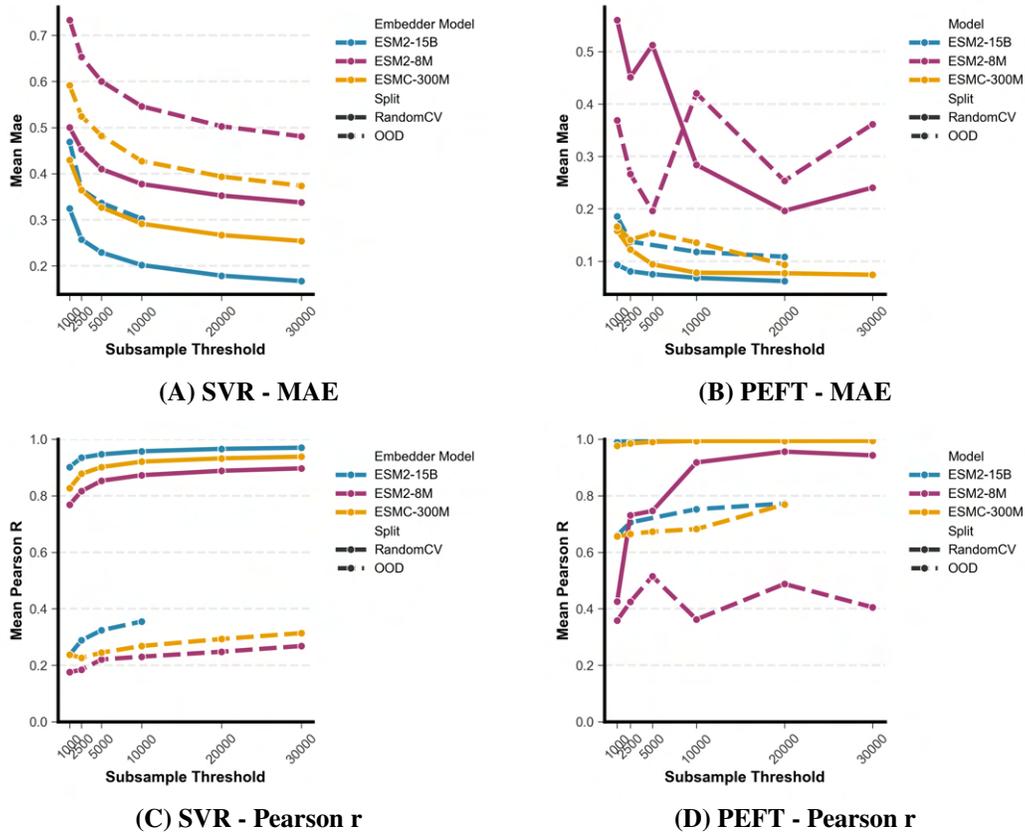


Figure 30: Sample size effects on predictive performance across model types and metrics. Panels show results for support vector regression (SVR, left) and PEFT models (right), evaluated using: (A,B) Mean Absolute Error (MAE), and (C,D) Pearson correlation as a function of increasing training set size, across RandomCV and OOD splits. Missing data points correspond to model collapse or cases exceeding a 48 hour training time limit.

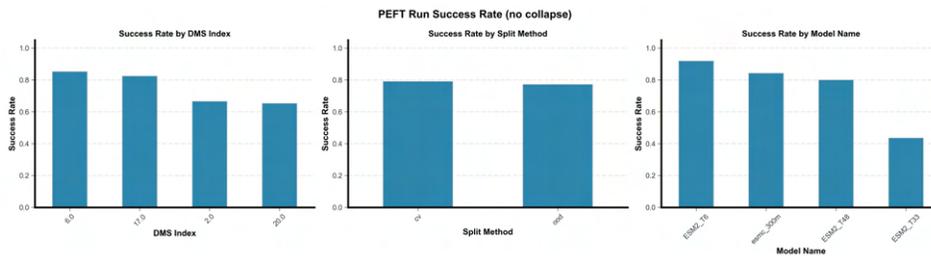


Figure 31: PEFT Training Success Rate

A.7 HIGH-THROUGHPUT SCREENING

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

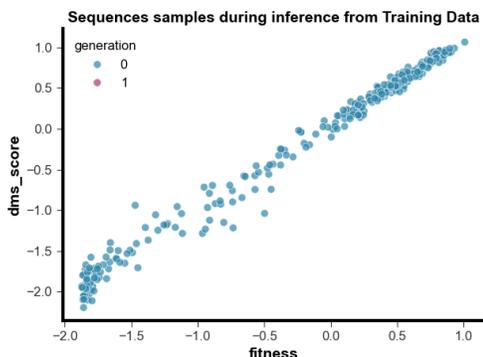


Figure 32: (A) Correlation of Inference matches of GB1\_IgG-Fc\_fitness\_1FCC over generations

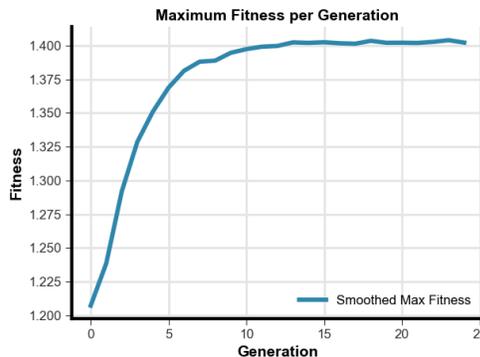


Figure 33: (B) Maximum fitness scores for GB1\_IgG-Fc\_fitness\_1FCC during inference across generations

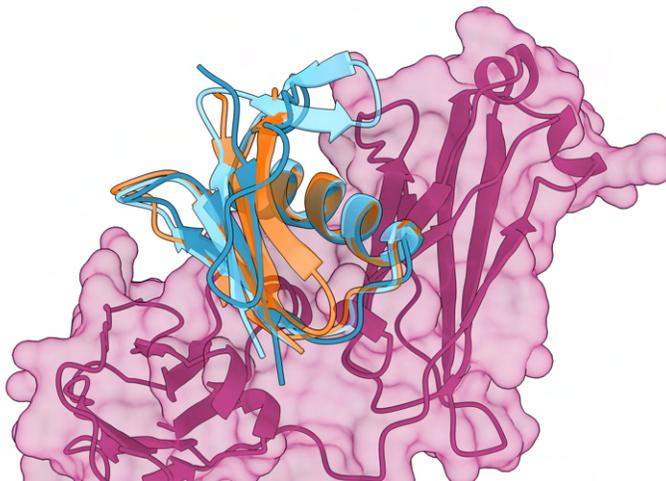


Figure 34: Overlaid binders from GB1\_IgG-Fc\_fitness\_1FCC. Showing the wild-type (orange), best from training (light blue) and best from inference (dark blue)



Figure 35: Sequence logo of the top 5 sequences of each generation for inference OOD split Trained ESMC-300M Ensemble on GB1\_IgG-Fc\_fitness\_1FCC.

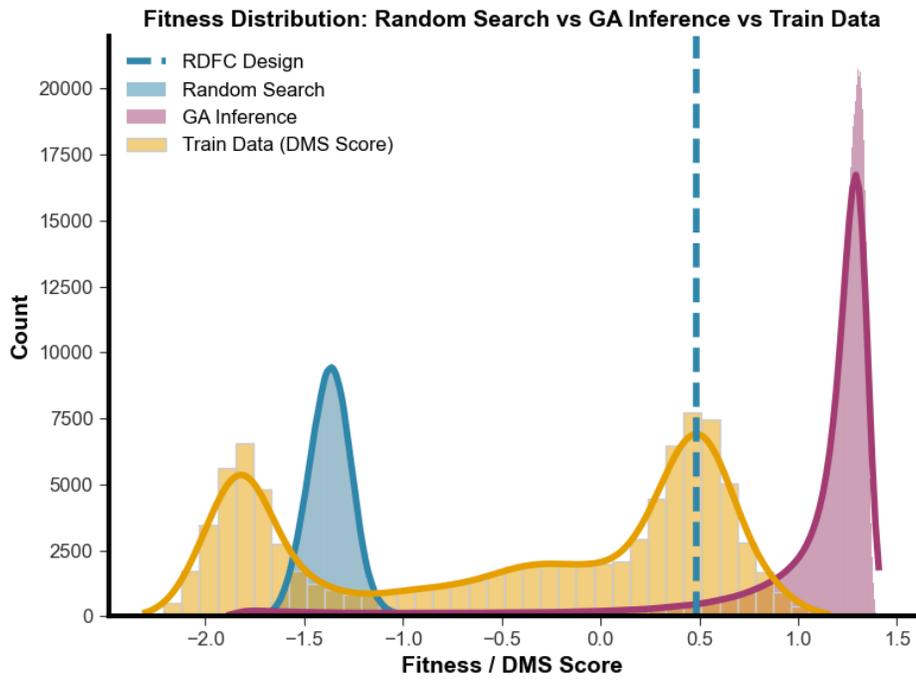


Figure 36: Fitness Distribution: Random Search vs GA Inference vs Train Data

### A.7.1 HYPERPARAMETER OPTIMIZATION

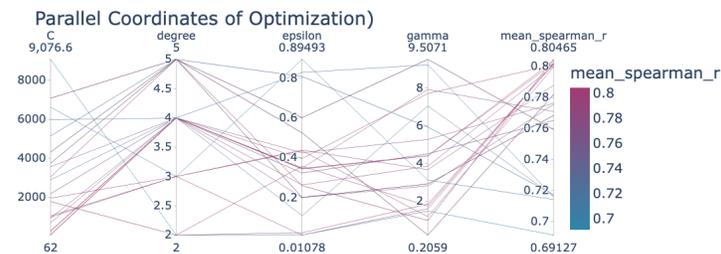


Figure 37: Parallel Coordinates of ACE2\_deltaKd (DMS index 8) Hyperparameter Optimization

1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

### Optimization History

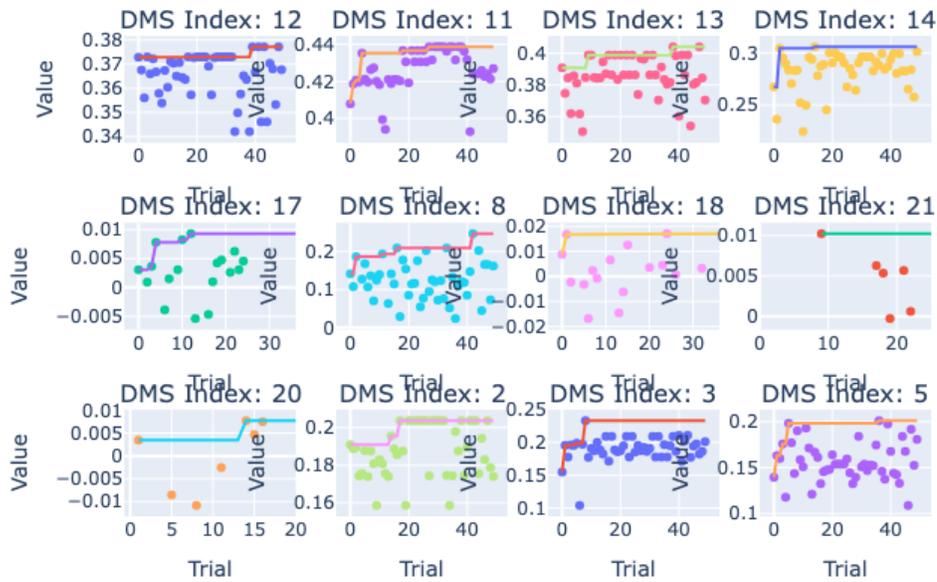


Figure 38: Optimization History for PEFT hyperparameters: rank, dropout

### Optimization History

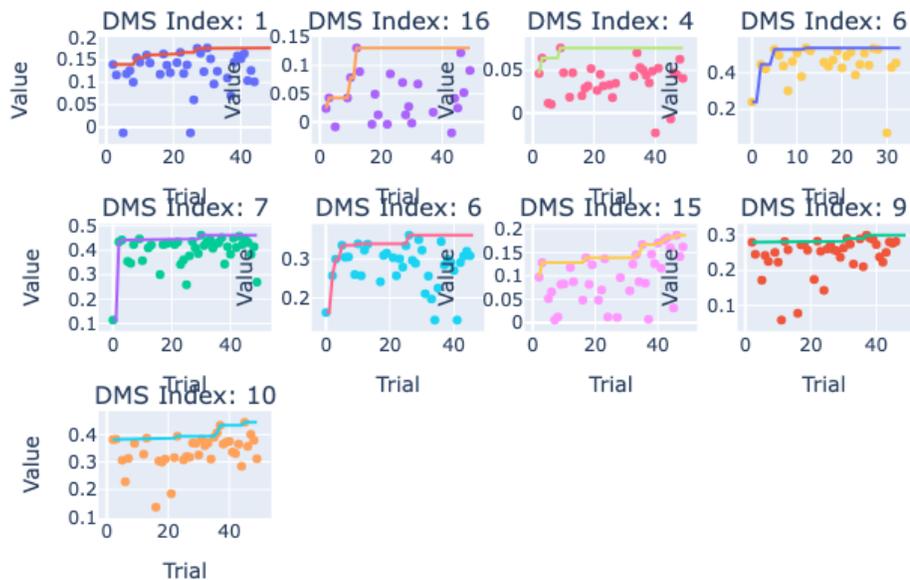


Figure 39: Optimization History for training hyperparameters: batch size, learning rate, weight decay, loss

## 1782 A.8 TABLES

1783

1784

1785

1786

1787

1788

1789

## A.8.1 TRAIN TEST SAMPLE COUNTS PER BENCHMARK AND SPLIT METHOD

1790

1791

1792

1793

1794

1795

1796

1797

1798

1799

1800

1801

1802

1803

1804

1805

Table 1: DMS Index 1

1806

1807

1808

1809

1810

1811

1812

1813

1814

1815

1816

1817

1818

1819

1820

1821

1822

1823

1824

1825

1826

1827

1828

1829

1830

1831

1832

1833

1834

1835

	RandomCV	ood	lomo	modulo
0	0: 23984/5997	0: 23957/6024	0: 11714/18267	0: 0/0
1	1: 23985/5996	1: 23905/6076	1: 18267/11714	1: 0/0
2	2: 23985/5996	2: 23827/6154	2: 2941/27040	2: 0/0
3	3: 23985/5996	3: 24129/5852	3: 18418/11563	3: 0/0
4	4: 23985/5996	4: 24106/5875	4: 5818/24163	4: 0/0
5			5: 5678/24303	
6			6: 4617/25364	
7			7: 3218/26763	
8			8: 4460/25521	
9			9: 5228/24753	
10			10: 2506/27475	
11			11: 2461/27520	
12			12: 5650/24331	
13			13: 3022/26959	
14			14: 2823/27158	
15			15: 3909/26072	
16			16: 3338/26643	
17			17: 2527/27454	
18			18: 2476/27505	
19			19: 12113/17868	
20			20: 2676/27305	
21			21: 10189/19792	
22			22: 6352/23629	
23			23: 4065/25916	
24			24: 5821/24160	
25			25: 2643/27338	
26			26: 15584/14397	
27			27: 2552/27429	
28			28: 4542/25439	
29			29: 7303/22678	

Table 2: DMS Index 2

	RandomCV	ood	lomo	modulo
1836				
1837				
1838				
1839	0	0: 36380/9096	0: 35957/9519	0: 7683/37793
1840	1	1: 36381/9095	1: 36576/8900	1: 8129/37347
1841	2	2: 36381/9095	2: 35713/9763	2: 5763/39713
1842	3	3: 36381/9095	3: 35349/10127	3: 8965/36511
1843	4	4: 36381/9095	4: 38309/7167	4: 14895/30581
1844	5		5: 8165/37311	
1845	6		6: 14304/31172	
1846	7		7: 3804/41672	
1847	8		8: 7291/38185	
1848	9		9: 11884/33592	
1849	10		10: 5703/39773	
1850	11		11: 13945/31531	
1851	12		12: 9666/35810	
1852	13		13: 12503/32973	
1853	14		14: 5211/40265	
1854	15		15: 12468/33008	
1855	16		16: 4500/40976	
1856	17		17: 13503/31973	
1857	18		18: 9772/35704	
1858	19		19: 6437/39039	
1859	20		20: 4314/41162	
1860	21		21: 20594/24882	
1861	22		22: 15865/29611	
1862	23		23: 16988/28488	
1863	24		24: 6485/38991	
1864	25		25: 12204/33272	
1865	26		26: 24730/20746	
1866	27		27: 14073/31403	
1867	28		28: 9618/35858	
1868	29		29: 15957/29519	
1869	30		30: 6119/39357	
1870	31		31: 12038/33438	
1871	32		32: 9030/36446	
1872	33		33: 4423/41053	
1873	34		34: 4634/40842	
1874	35		35: 3823/41653	
1875				
1876				
1877				
1878				
1879				
1880				
1881				
1882				
1883				
1884				
1885				
1886				
1887				
1888				
1889				

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

Table 3: DMS Index 3

	RandomCV	ood	lomo	modulo
0	0: 4466/1117	0: 4446/1137	0: 1023/4560	0: 0/0
1	1: 4466/1117	1: 4481/1102	1: 829/4754	1: 0/0
2	2: 4466/1117	2: 4471/1112	2: 1395/4188	2: 0/0
3	3: 4467/1116	3: 4467/1116	3: 550/5033	3: 0/0
4	4: 4467/1116	4: 4467/1116	4: 1313/4270	4: 0/0
5			5: 941/4642	
6			6: 925/4658	
7			7: 452/5131	
8			8: 1288/4295	
9			9: 828/4755	
10			10: 1155/4428	
11			11: 904/4679	
12			12: 939/4644	
13			13: 644/4939	
14			14: 1055/4528	
15			15: 485/5098	
16			16: 611/4972	
17			17: 522/5061	
18			18: 3724/1859	
19			19: 474/5109	
20			20: 789/4794	
21			21: 1429/4154	
22			22: 3547/2036	
23			23: 4920/663	
24			24: 3371/2212	
25			25: 890/4693	
26			26: 835/4748	

Table 4: DMS Index 4

	RandomCV	ood	lomo	modulo
0	0: 4468/1117	0: 4468/1117		0: 0/0
1	1: 4468/1117	1: 4468/1117		1: 0/0
2	2: 4468/1117	2: 4468/1117		2: 0/0
3	3: 4468/1117	3: 4468/1117		3: 0/0
4	4: 4468/1117	4: 4468/1117		4: 0/0

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

Table 5: DMS Index 5

	RandomCV	ood	lomo	modulo
0	0: 14725/3682	0: 14682/3725	0: 14590/3817	0: 0/0
1	1: 14725/3682	1: 14664/3743	1: 16000/2407	1: 0/0
2	2: 14726/3681	2: 14831/3576	2: 15935/2472	2: 0/0
3	3: 14726/3681	3: 14725/3682	3: 16577/1830	3: 0/0
4	4: 14726/3681	4: 14726/3681	4: 16142/2265	4: 0/0
5			5: 16838/1569	
6			6: 16905/1502	
7			7: 15879/2528	
8			8: 16500/1907	
9			9: 16292/2115	
10			10: 14804/3603	
11			11: 15780/2627	
12			12: 16820/1587	
13			13: 16019/2388	
14			14: 16691/1716	

Table 6: DMS Index 6

	RandomCV	ood	lomo	modulo
0	0: 74312/18579	0: 74311/18580	0: 84805/8086	0: 0/0
1	1: 74313/18578	1: 74314/18577		1: 0/0
2	2: 74313/18578	2: 74313/18578		2: 0/0
3	3: 74313/18578	3: 74313/18578		3: 0/0
4	4: 74313/18578	4: 74313/18578		4: 0/0

Table 7: DMS Index 7

	RandomCV	ood	lomo	modulo
0	0: 17740/4436	0: 17740/4436	0: 2024/20152	0: 0/0
1	1: 17741/4435	1: 17741/4435	1: 3642/18534	1: 0/0
2	2: 17741/4435	2: 17741/4435	2: 4194/17982	2: 0/0
3	3: 17741/4435	3: 17741/4435	3: 2416/19760	3: 0/0
4	4: 17741/4435	4: 17741/4435	4: 1853/20323	4: 0/0
5			5: 2197/19979	
6			6: 2353/19823	
7			7: 1820/20356	
8			8: 2934/19242	
9			9: 2390/19786	
10			10: 2141/20035	
11			11: 2622/19554	
12			12: 2508/19668	
13			13: 2326/19850	
14			14: 2128/20048	
15			15: 1890/20286	
16			16: 2328/19848	
17			17: 2186/19990	

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

Table 8: DMS Index 8

	RandomCV	ood	lomo	modulo
0	0: 17497/4375	0: 16465/5407		0: 0/0
1	1: 17497/4375	1: 18526/3346		1: 0/0
2	2: 17498/4374	2: 17466/4406		2: 0/0
3	3: 17498/4374	3: 17503/4369		3: 0/0
4	4: 17498/4374	4: 17528/4344		4: 0/0

Table 9: DMS Index 9

	RandomCV	ood	lomo	modulo
0	0: 15626/3907	0: 15626/3907	0: 17272/2261	0: 0/0
1	1: 15626/3907	1: 15627/3906	1: 17597/1936	1: 0/0
2	2: 15626/3907	2: 15626/3907	2: 1667/17866	2: 0/0
3	3: 15627/3906	3: 15627/3906		3: 0/0
4	4: 15627/3906	4: 15626/3907		4: 0/0

Table 10: DMS Index 10

	RandomCV	ood	lomo	modulo
0	0: 15362/3841	0: 15362/3841	0: 17270/1933	0: 0/0
1	1: 15362/3841	1: 15363/3840	1: 17292/1911	1: 0/0
2	2: 15362/3841	2: 15362/3841	2: 1637/17566	2: 0/0
3	3: 15363/3840	3: 15363/3840		3: 0/0
4	4: 15363/3840	4: 15362/3841		4: 0/0

Table 11: DMS Index 11

	RandomCV	ood	lomo	modulo
0	0: 10141/2536	0: 10141/2536	0: 11617/1060	0: 0/0
1	1: 10141/2536	1: 10142/2535	1: 11630/1047	1: 0/0
2	2: 10142/2535	2: 10142/2535	2: 10316/2361	2: 0/0
3	3: 10142/2535	3: 10142/2535	3: 11652/1025	3: 0/0
4	4: 10142/2535	4: 10141/2536	4: 11646/1031	4: 0/0
5			5: 11510/1167	
6			6: 11508/1169	

Table 12: DMS Index 12

	RandomCV	ood	lomo	modulo
0	0: 18529/4633	0: 18529/4633	0: 20564/2598	0: 0/0
1	1: 18529/4633	1: 18530/4632	1: 21008/2154	1: 0/0
2	2: 18530/4632	2: 18530/4632	2: 1870/21292	2: 0/0
3	3: 18530/4632	3: 18530/4632		3: 0/0
4	4: 18530/4632	4: 18529/4633		4: 0/0

2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105

Table 13: DMS Index 13

	RandomCV	ood	lomo	modulo
0	0: 16272/4069	0: 16272/4069	0: 17829/2512	0: 0/0
1	1: 16273/4068	1: 16273/4068	1: 18390/1951	1: 0/0
2	2: 16273/4068	2: 16273/4068	2: 1701/18640	2: 0/0
3	3: 16273/4068	3: 16273/4068		3: 0/0
4	4: 16273/4068	4: 16273/4068		4: 0/0

Table 14: DMS Index 14

	RandomCV	ood	lomo	modulo
0	0: 15540/3885	0: 15540/3885	0: 17164/2261	0: 0/0
1	1: 15540/3885	1: 15540/3885	1: 17532/1893	1: 0/0
2	2: 15540/3885	2: 15540/3885	2: 1624/17801	2: 0/0
3	3: 15540/3885	3: 15540/3885		3: 0/0
4	4: 15540/3885	4: 15540/3885		4: 0/0

Table 15: DMS Index 15

	RandomCV	ood	lomo	modulo
0	0: 2675/669	0: 2675/669		0: 0/0
1	1: 2675/669	1: 2675/669		1: 0/0
2	2: 2675/669	2: 2676/668		2: 0/0
3	3: 2675/669	3: 2675/669		3: 0/0
4	4: 2676/668	4: 2675/669		4: 0/0

Table 16: DMS Index 16

	RandomCV	ood	lomo	modulo
0	0: 3108/778	0: 3108/778	0: 2460/1426	0: 0/0
1	1: 3109/777	1: 3109/777	1: 2464/1422	1: 0/0
2	2: 3109/777	2: 3105/781	2: 2443/1443	2: 0/0
3	3: 3109/777	3: 3105/781		3: 0/0
4	4: 3109/777	4: 3117/769		4: 0/0

Table 17: DMS Index 17

	RandomCV	ood	lomo	modulo
2106				
2107				
2108				
2109	0	0: 26052/6513	0: 16290/16275	0: 0/0
2110	1	1: 26052/6513	1: 16275/16290	1: 0/0
2111	2	2: 26052/6513	2: 16289/16276	2: 0/0
2112	3	3: 26052/6513	3: 16276/16289	3: 0/0
2113	4	4: 26052/6513	4: 16291/16274	4: 0/0
2114	5		5: 16274/16291	
2115	6		6: 16349/16216	
2116	7		7: 16216/16349	
2117	8		8: 16380/16185	
2118	9		9: 16185/16380	
2119	10		10: 16294/16271	
2120	11		11: 16271/16294	
2121	12		12: 16307/16258	
2122	13		13: 16258/16307	
2123	14		14: 16344/16221	
2124	15		15: 16221/16344	
2125	16		16: 16304/16261	
2126	17		17: 16261/16304	
2127	18		18: 16249/16316	
2128	19		19: 16316/16249	
2129	20		20: 16282/16283	
2130	21		21: 16283/16282	
2131	22		22: 16374/16191	
2132	23		23: 16191/16374	
2133	24		24: 16181/16384	
2134	25		25: 16384/16181	
2135	26		26: 16181/16384	
2136	27		27: 16384/16181	
2137	28		28: 16181/16384	
2138	29		29: 16384/16181	
2139				
2140				
2141				
2142				
2143				
2144				
2145				
2146				
2147				
2148				
2149				
2150				
2151				
2152				
2153				
2154				
2155				
2156				
2157				
2158				
2159				

Table 18: DMS Index 18

	RandomCV	ood	lomo	modulo
2160				
2161				
2162				
2163	0	0: 13208/3303	0: 8275/8236	0: 0/0
2164	1	1: 13209/3302	1: 8236/8275	1: 0/0
2165	2	2: 13209/3302	2: 8233/8278	2: 0/0
2166	3	3: 13209/3302	3: 8278/8233	3: 0/0
2167	4	4: 13209/3302	4: 8261/8250	4: 0/0
2168	5		5: 8250/8261	
2169	6		6: 8197/8314	
2170	7		7: 8314/8197	
2171	8		8: 8264/8247	
2172	9		9: 8247/8264	
2173	10		10: 8255/8256	
2174	11		11: 8256/8255	
2175	12		12: 8237/8274	
2176	13		13: 8274/8237	
2177	14		14: 8264/8247	
2178	15		15: 8247/8264	
2179	16		16: 8211/8300	
2180	17		17: 8300/8211	
2181	18		18: 8297/8214	
2182	19		19: 8214/8297	
2183	20		20: 8287/8224	
2184	21		21: 8224/8287	
2185	22		22: 8316/8195	
2186	23		23: 8195/8316	
2187	24		24: 8305/8206	
2188	25		25: 8206/8305	
2189	26		26: 8227/8284	
2190	27		27: 8284/8227	
2191				
2192				
2193				
2194				
2195				
2196				
2197				
2198				
2199				
2200				
2201				
2202				
2203				
2204				
2205				
2206				
2207				
2208				
2209				
2210				
2211				
2212				
2213				

Table 19: DMS Index 19

	RandomCV	ood	lomo	modulo	
2214					
2215					
2216					
2217					
2218	0	0: 15893/3974	0: 15893/3974	0: 10312/9555	0: 0/0
2219	1	1: 15893/3974	1: 15894/3973	1: 9555/10312	1: 0/0
2220	2	2: 15894/3973	2: 15894/3973	2: 9792/10075	2: 0/0
2221	3	3: 15894/3973	3: 15894/3973	3: 10075/9792	3: 0/0
2222	4	4: 15894/3973	4: 15893/3974	4: 10048/9819	4: 0/0
2223	5			5: 9819/10048	
2224	6			6: 8359/11508	
2225	7			7: 11508/8359	
2226	8			8: 9499/10368	
2227	9			9: 10368/9499	
2228	10			10: 10020/9847	
2229	11			11: 9847/10020	
2230	12			12: 10847/9020	
2231	13			13: 9020/10847	
2232	14			14: 9759/10108	
2233	15			15: 10108/9759	
2234	16			16: 9741/10126	
2235	17			17: 10126/9741	
2236	18			18: 6594/13273	
2237	19			19: 13273/6594	
2238	20			20: 14525/5342	
2239	21			21: 5342/14525	
2240	22			22: 9904/9963	
2241	23			23: 9963/9904	
2242	24			24: 10118/9749	
2243	25			25: 9749/10118	
2244	26			26: 11174/8693	
2245	27			27: 8693/11174	
2246	28			28: 10074/9793	
2247	29			29: 9793/10074	
2248					
2249					
2250					
2251					
2252					
2253					
2254					
2255					
2256					
2257					
2258					
2259					
2260					
2261					
2262					
2263					
2264					
2265					
2266					
2267					

Table 20: DMS Index 20

	RandomCV	ood	lomo	modulo
2268				
2269				
2270				
2271	0	0: 26214/6554	0: 16384/16384	0: 0/0
2272	1	1: 26214/6554	1: 16384/16384	1: 0/0
2273	2	2: 26214/6554	2: 16384/16384	2: 0/0
2274	3	3: 26215/6553	3: 16384/16384	3: 0/0
2275	4	4: 26215/6553	4: 16384/16384	4: 0/0
2276	5		5: 16384/16384	
2277	6		6: 16384/16384	
2278	7		7: 16384/16384	
2279	8		8: 16384/16384	
2280	9		9: 16384/16384	
2281	10		10: 16384/16384	
2282	11		11: 16384/16384	
2283	12		12: 16384/16384	
2284	13		13: 16384/16384	
2285	14		14: 16384/16384	
2286	15		15: 16384/16384	
2287	16		16: 16384/16384	
2288	17		17: 16384/16384	
2289	18		18: 16384/16384	
2290	19		19: 16384/16384	
2291	20		20: 16384/16384	
2292	21		21: 16384/16384	
2293	22		22: 16384/16384	
2294	23		23: 16384/16384	
2295	24		24: 16384/16384	
2296	25		25: 16384/16384	
2297	26		26: 16384/16384	
2298	27		27: 16384/16384	
2299	28		28: 16384/16384	
2300	29		29: 16384/16384	
2301				
2302				
2303				
2304				
2305				
2306				
2307				
2308				
2309				
2310				
2311				
2312				
2313				
2314				
2315				
2316				
2317				
2318				
2319				
2320				
2321				

Table 21: DMS Index 21

	RandomCV	ood	lomo	modulo
2322				
2323				
2324				
2325	0	0: 18948/4738	0: 12051/11635	0: 0/0
2326	1	1: 18949/4737	1: 11635/12051	1: 0/0
2327	2	2: 18949/4737	2: 11655/12031	2: 0/0
2328	3	3: 18949/4737	3: 12031/11655	3: 0/0
2329	4	4: 18949/4737	4: 11929/11757	4: 0/0
2330	5		5: 11757/11929	
2331	6		6: 11183/12503	
2332	7		7: 12503/11183	
2333	8		8: 11760/11926	
2334	9		9: 11926/11760	
2335	10		10: 8626/15060	
2336	11		11: 15060/8626	
2337	12		12: 16376/7310	
2338	13		13: 7310/16376	
2339	14		14: 11849/11837	
2340	15		15: 11837/11849	
2341	16		16: 11834/11852	
2342	17		17: 11852/11834	
2343	18		18: 12211/11475	
2344	19		19: 11475/12211	
2345	20		20: 11751/11935	
2346	21		21: 11935/11751	
2347	22		22: 11922/11764	
2348	23		23: 11764/11922	
2349	24		24: 12180/11506	
2350	25		25: 11506/12180	
2351	26		26: 11621/12065	
2352	27		27: 12065/11621	
2353	28		28: 11941/11745	
2354	29		29: 11745/11941	
2355				
2356				
2357				
2358				
2359				
2360				
2361				
2362				
2363				
2364				
2365				
2366				
2367				
2368				
2369				
2370				
2371				
2372				
2373				
2374				
2375				

## A.8.2 MAPPING FROM DMS INTEGER TO DMS ID

Table 22: Mapping from DMS integer to DMS ID

DMS Index	DMS_id
1	5A12_VEGF_fitness_4ZFF
2	Z-domain_ZSPA-1_LL1_fitness_1LP1
3	Z-domain_ZSPA-1_LL2_fitness_1LP1
4	CXCR4_CXCL12_enrich_8U4O
5	hYAP65_peptide_FunctionalScore_1JMQ
6	GB1_IgG-Fc_fitness_1FCC
7	GB1_IgG-Fc_fitness_1FCC_2016
8	SARS2-RBD_ACE2_deltaKd_6M0J
9	KRAS_DARPinK27_norfitness_5O2S
10	KRAS_PICK3CG-RBD_norfitness_1HE8
11	KRAS_RAF1_norfitness_6VJJ
12	KRAS_RAF1-RBD_norfitness_6VJJ
13	KRAS_RALGDS-RBD_norfitness_1LFD
14	KRAS_SOS1_norfitness_8BE4
15	HLA-A2_TAPBPR_meanscore_5WER
16	CD19_FMC63_Fitness_7URV
17	SARS2-RBD_ACE2_HUMAN_7WPB
18	SARS2-RBD_LY-CoV016_7C01
19	SARS2-RBD_LY-CoV555_7KMG
20	SARS2-RBD_S309_7XCO
21	SARS2-RBD_REGN10987_9LYP

## A.8.3 MODEL TRAINABLE PARAMETERS

Table 23: Total and Trainable Number of Parameters

Model	Trainable	Total	Source
ESM2-T6	31.4K	8M	Lin et al. (2022)
ESM2-T12	90.3K	35M	Lin et al. (2022)
ESM2-T30	293K	150M	Lin et al. (2022)
ESM2-T33	643K	650M	Lin et al. (2022)
ESM2-T36	1.4M	3B	Lin et al. (2022)
ESM2-T48	3.7M	15B	Lin et al. (2022)
ESMC-300M	929K	300M	github
ESMC-600M	1.3M	600M	github

## A.8.4 ABSOLUTE MEAN DIFFERENCE BY MUTATION FOR COMBINATORIAL LIBRARIES

Table 24: Absolute mean differences of labels by mutation for each Combinatorial Dataset

Mutation Index	20 abs_mean_diff	17 abs_mean_diff	21 abs_mean_diff	19 abs_mean_diff	18 abs_mean_diff
1	0.3270	0.0021	0.0424	0.0849	0.0790
2	0.1726	0.0340	0.0091	0.0329	0.0145
3	0.2176	0.0084	0.0336	0.0526	0.0647
4	0.1996	0.0959	0.2034	0.3542	0.3161
5	0.0650	0.3046	0.0703	0.0363	1.3200
6	0.0131	0.1273	0.1950	0.0436	0.0225
7	0.0325	0.2027	0.6392	0.0360	0.0227
8	0.0029	0.2756	0.0109	0.0099	0.0405
9	0.0093	0.0234	0.0035	0.0791	0.0773
10	0.0361	0.0263	0.1407	1.0990	0.1698
11	0.0246	0.1290	0.0112	1.7586	0.9759
12	0.0227	0.2204	0.1706	0.0290	0.1994
13	0.0851	0.3436	0.0231	0.1533	0.1746
14	0.0132	1.0585	0.1223	0.2738	0.2057
15	0.0134	0.3478	0.0121	0.0986	0.0469

## A.8.5 RANDOM FOREST RESULTS

Table 25: Random Forest Results for ESMC-300M Embeddings on OOD splits

DMS Index	kendall tau	pearson r	spearman r	mse	mae	topk 10
18	0.434009	0.597546	0.637233	0.253584	0.409949	0.159903
8	0.083003	0.090745	0.128024	3.742846	1.668905	0.182442
17	0.607948	0.812376	0.833095	0.190022	0.347576	0.241093
14	0.416928	0.620189	0.637868	0.129266	0.288854	0.032956
13	0.493056	0.749092	0.721773	0.093235	0.253618	0.128810
11	0.453245	0.727540	0.689020	0.126496	0.284785	0.142068
12	0.509364	0.756736	0.741292	0.067297	0.211846	0.180484
10	0.444506	0.670757	0.669320	0.139110	0.303964	0.080729
9	0.416221	0.617737	0.636157	0.129906	0.290531	0.036354
15	NaN	NaN	NaN	1.436910	0.928248	0.000000
7	0.503302	0.661358	0.732661	0.628985	0.575291	0.152007
6	0.477235	0.737131	0.689074	0.485125	0.569585	0.070083
4	NaN	NaN	NaN	2.983544	1.364413	0.000000
16	0.318336	0.681855	0.576211	13.264859	2.832486	0.000000
1	0.340525	0.901817	0.578464	0.187763	0.307552	0.000000
4	NaN	NaN	NaN	2.966439	1.358867	0.000000
16	0.319696	0.683601	0.576988	13.208873	2.823042	0.000000
1	0.338679	0.901859	0.576989	0.187721	0.307030	0.000000

Table 26: Random Forest Results for ESM2-8M Embeddings on OOD splits

DMS Index	kendall tau	pearson r	spearman r	mse	mae	topk 10
19	NaN	NaN	NaN	NaN	NaN	NaN
18	0.397687	0.560701	0.589294	0.265554	0.418347	0.117505
8	NaN	NaN	NaN	4.145915	1.783551	0.122542
17	0.504306	0.683923	0.741008	0.272118	0.421938	0.141892
14	0.325087	0.481006	0.514187	0.161161	0.328247	0.008754
13	0.410620	0.640304	0.616979	0.125702	0.299519	0.065880
11	0.345055	0.540886	0.554142	0.184435	0.342490	0.037885
12	0.425853	0.643406	0.636441	0.092589	0.255073	0.101468
10	0.345122	0.544045	0.541973	0.177010	0.349289	0.055208
9	0.329740	0.483978	0.517242	0.160598	0.328371	0.020993
15	NaN	NaN	NaN	1.449088	0.930009	0.000000
7	0.430133	0.530310	0.647429	0.764871	0.653123	0.059540
6	0.449369	0.693868	0.651968	0.552429	0.616058	0.098396
4	NaN	NaN	NaN	2.914910	1.341311	0.000000
16	0.307378	0.668956	0.560441	13.698226	2.890746	0.000000
1	0.353628	0.861195	0.582884	0.256724	0.360615	0.000000
6	0.450873	0.695145	0.653538	0.550782	0.615207	0.104209
4	NaN	NaN	NaN	2.911214	1.339890	0.003584
16	0.307471	0.670012	0.560837	13.662221	2.884809	0.000000
1	0.353752	0.862003	0.583734	0.255407	0.359599	0.000000

## A.8.6 SVR-PROTEINMPNN RESULTS

Table 27: SVR-MPNN Results on OOD splits

DMS Index	spearman r	pearson r	mae	rmse	topk 10	kendall tau
16	0.077831	0.109733	6.575633	6.671940	0.152685	0.052634
4	0.025040	0.018464	1.347288	1.400675	0.102703	0.016541
6	0.116064	0.107340	1.253058	1.277679	0.134181	0.077876
7	0.105954	0.108899	0.899254	0.930745	0.162077	0.071175
15	0.039693	0.038383	0.995695	1.046435	0.106061	0.026854
9	0.074369	0.081582	0.561248	0.570047	0.115897	0.049625
10	0.008700	0.034675	0.625334	0.634298	0.093229	0.005916
12	0.021951	0.034896	0.478465	0.484774	0.092873	0.014534
11	0.078024	0.085800	0.573439	0.582957	0.118577	0.052013
13	0.028497	0.037549	0.588771	0.597637	0.105419	0.018843
14	0.059904	0.073836	0.561718	0.570508	0.123711	0.039887
17	0.153440	0.166940	0.774771	0.789607	0.160369	0.103226
8	0.046785	0.046397	2.062201	2.119056	0.113516	0.031241
18	0.222224	0.158925	0.774861	0.793422	0.167879	0.148904
21	0.146686	0.166995	0.489976	0.508612	0.158140	0.099702
2	NaN	NaN	0.096960	0.125291	0.114767	NaN
3	0.136143	0.129061	0.166844	0.223637	0.111668	0.091136
5	0.067119	0.069952	1.052114	1.125921	0.119786	0.044952

## A.8.7 BEST TRIALS FOR PEFT RANK AND DROPOUT OPTIMIZATION

Table 28: Best Trials Parameters PEFT

	DMS Index	spearman r	rank	dropout
	0	12	0.377184	2 0.200000
	1	11	0.438652	32 0.100000
	2	13	0.404050	32 0.300000
	3	14	0.305758	1 0.300000
	4	17	0.009290	4 0.200000
	5	8	0.247150	1 0.200000
	6	18	0.017090	1 0.400000
	7	21	0.010211	32 0.300000
	8	20	0.007740	1 0.300000
	9	2	0.203734	2 0.300000
	10	3	0.233557	1 0.100000
	11	5	0.201198	2 0.200000

## A.8.8 BEST TRIALS FOR TRAINING HYPERPARAMETER OPTIMIZATION

Table 29: Best Trials Parameters Training

DMS Index	spearman r	batch size	learning rate	weight decay	loss
	1	0.176577	128	0.000341	0.059111 mse
	16	0.130845	16	0.000253	0.049775 huber
	4	0.076111	128	0.000230	0.007790 l1
	6	0.536756	16	0.000253	0.049775 huber
	7	0.462760	8	0.000109	0.069547 l1
	6	0.362466	16	0.000554	0.074595 huber
	15	0.187319	4	0.000072	0.078060 mse
	9	0.300291	4	0.000107	0.057737 mse
	10	0.444399	32	0.000337	0.046022 huber

## A.9 PREDICTOR EXAMPLE

**Minimal predictor example implementing all abstract methods:**

Listing 1: Minimal example: a predictor returning the mean of training labels

```

2577 import numpy as np
2578 from haipr.predictor import BasePredictor
2579
2580 class MyPredictor(BasePredictor):
2581     def __init__(self):
2582         super(). __init__()
2583         self.mean_value = None
2584
2584     def setup_model(self, data, cfg):
2585         self.data = data
2586         self.cfg = cfg
2587
2588     def fit_model(self, dataset, train_indices, val_indices=None):
2589         # Save mean of training labels
2590         labels = np.array(dataset.get_labels())
2591         self.mean_value = float(np.mean(labels[train_indices]))
2591         mean_array = np.full(len(train_indices), self.mean_value)

```

```

2592         metrics = compute_regression_metrics(labels[train_indices], mean_array)
2593         predictions = { "indices": train_indices,
2594                       "predictions": mean_array,
2595                       "true_values": labels[train_indices]}
2596         return {"metrics": metrics, "predictions": predictions}
2597
2598     def predict_sequences(self, sequences, params=None):
2599         # Always return self.mean_value for each sequence
2600         return np.full(len(sequences), self.mean_value)
2601
2602     def save_model(self, save_dir):
2603         # Dummy implementation to satisfy interface
2604         return None
2605
2606     def prepare_training_features(self, dataset, indices):
2607         pass
2608
2609     def prepare_batch_features(self, batch_items):
2610         pass

```

#### YAML configuration for the predictor in conf/models/my\_predictor.yaml:

Listing 2: Configuration example for MyPredictor

```

2614 defaults:
2615   - _self_
2616   _target_: path.to.your.MyPredictor

```

Running the new predictor is as simple as running the following command:

Listing 3: Running the new predictor

```

2620 python -m haipr.haipr stages=train model=my_predictor

```

#### A.10 GENERATOR EXAMPLE

##### Minimal sequence generator example implementing all abstract methods from BaseSequenceGenerator:

Listing 4: Minimal example of a custom sequence generator

```

2628 from haipr.sequence_generators.base_generator import BaseSequenceGenerator
2629
2630 class MySequenceGenerator(BaseSequenceGenerator):
2631     def __init__(self, num_sequences=100):
2632         self.data = None
2633         self.num_sequences = num_sequences
2634         self.alphabet_per_position = None
2635         self.fitness_callback = None
2636         self.best_sequence = None
2637         self.best_fitness = None
2638         self.all_sequences = []
2639         self.all_fitnesses = []
2640         self.metrics_logger = None
2641         self.new_run_callback = None
2642
2643     def setup_generator(self, data, alphabet_per_position, fitness_callback):
2644         self.data = data
2645         self.alphabet_per_position = alphabet_per_position
2646         self.fitness_callback = fitness_callback

```

```

2646     def run_generator(self):
2647         for i in range(self.num_sequences):
2648             random_sequence = ''.join([random.choice(self.alphabet_per_position[i])
2649             self.all_sequences.append(random_sequence)
2650             fitness = self.fitness_callback(random_sequence)
2651             self.all_fitnesses.append(fitness)
2652             self.metrics_logger({"fitness": fitness}, step=i)
2653             if fitness > self.best_fitness:
2654                 self.best_sequence = random_sequence
2655                 self.best_fitness = fitness
2656         return self.all_sequences, self.all_fitnesses
2657
2658     def get_best_solution(self):
2659         return self.best_sequence, self.best_fitness
2660
2661     def shutdown(self):
2662         # No-op for this minimal example
2663         pass
2664
2665     def set_metrics_logger(self, logger_func):
2666         self.metrics_logger = logger_func

```

#### YAML configuration for the predictor in conf/models/my\_generator.yaml:

Listing 5: Configuration example for MySequenceGenerator

```

2669 defaults:
2670   - _self_
2671   _target_: path.to.your.MySequenceGenerator
2672   num_sequences: 100

```

Running the new predictor is as simple as running the following command:

Listing 6: Running the new predictor

```
python -m haipr.haipr stages=inference generator=my_generator
```

## A.11 USE OF LLMs

We used LLMs to aid in preventing repetitive words and optimize sentence structure as well as language.

2683  
2684  
2685  
2686  
2687  
2688  
2689  
2690  
2691  
2692  
2693  
2694  
2695  
2696  
2697  
2698  
2699