
Sparse Autoencoders in Protein Engineering Campaigns: Steering and Model Diffing

Gerard Boxó¹ Filippo Stocco^{2,1} Noelia Ferruz^{1,2}

Abstract

Protein Language Models (pLM) have proven versatile tools in protein design, but their internal workings remain difficult to interpret. Here, we implement a mechanistic interpretability framework and apply it in two scenarios. First, by training sparse autoencoders (SAEs) on the model activations, we identify and annotate features relevant to enzyme variant activity through a two-stage process involving candidate selection and causal intervention. During sequence generation, we steer the model by clamping or ablating key SAE features, which increases the predicted enzyme activity. Additionally, we implement an intervention strategy: *MSA-steering*, which projects SAE latents in the multiple sequence alignment dimensionality of our case study enzyme. Second, we compare pLM checkpoints before and after three rounds of Reinforcement Learning (RL) by examining sequence regions with high divergence of per-token log-likelihood, detecting the residues that most align with higher predicted affinities. Overall, we present a strategy to apply SAE for protein engineering.

1. Introduction

End-to-end differentiable models are complex nonlinear functions $f : X \rightarrow Y$ that map an input space X to an output space Y . These mapping functions are essentially black boxes, making it difficult to explain *how* and *why* a model ends up making a particular decision. Protein language models (pLMs), are no exception, but despite their hermetic nature, pLMs must have nevertheless learned some

complex sequence-to-function relationships, as evidenced by their versatility and state-of-the-art performance in tasks ranging from protein folding (Lin et al., 2023a) to protein design (Yang et al., 2024; Madani et al., 2023; Bhatnagar et al., 2025), including distant yet catalytically efficient enzymes (Munsamy et al., 2022; Madani et al., 2023; Johnson et al., 2023; Parsan et al., 2025).

Mechanistic interpretability aims to provide a detailed analysis of the mechanisms underlying the predictions of deep learning models. Sparse Autoencoders (SAE) in particular have recently emerged as a relevant tool to extract interpretable features, for the study of internal circuits from LLMs (Marks et al., 2024). In the field of protein research, we are witnessing applications for pLMs with promising outcomes, especially for the understanding of encoder-only pLMs (Parsan et al., 2025; Simon & Zou, 2024; Adams et al., 2025; Garcia & Ansuini, 2025).

SAE models consist of an encoder-decoder architecture that learns to produce intermediate activations of higher dimensionality (1), incentivized to be sparse through the training process. In particular, the encoder transforms an input x into an intermediate vector through a function f , ensuring the activations are sparse (i.e., present few non-zero features) by applying a *BatchTopK* activation that retains the $k \times n$ largest entries of the SAE latent within each batch, zeroing out all the others (Bussmann et al., 2024) (Eq. 1). The decoder learns to reconstruct the activations x as output (Eq. 2), by applying a training loss that is formulated to both reconstruct the model activation by the mean square error of the vector x and \hat{x} (Eq. 3) with the auxiliary loss that ensures sparsity (Eq. 4):

$$\mathbf{f}(\mathbf{x}) = \text{BatchTopK}(\mathbf{W}_{\text{enc}}\mathbf{x} + \mathbf{b}_{\text{enc}}) \quad (1)$$

$$\hat{\mathbf{x}} = \mathbf{W}_{\text{dec}}\mathbf{f} + \mathbf{b}_{\text{dec}} \quad (2)$$

$$\mathcal{L}(\mathbf{x}) = \text{MSE}(\mathbf{x}, \hat{\mathbf{x}}) + \mathcal{L}_{\text{aux}} \quad (3)$$

$$\mathcal{L}_{\text{aux}} = \text{MSE}(\mathbf{e}, \hat{\mathbf{e}}) \quad (4)$$

In this work, we investigate the potential of SAEs in the context of decoder-only pLMs. We explore their appli-

¹Centre for Genomic Regulation, The Barcelona Institute of Science and Technology, Barcelona, Spain ²Department of Medicine and Life Sciences, Universitat Pompeu Fabra, Barcelona, Spain. Correspondence to: Gerard Boxó; Noelia Ferruz <gerard.boxo@crg.eu; noelia.ferruz@crg.eu>.

Proceedings of the Workshop on Generative AI for Biology at the 42nd International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

cation for interventions during inference (steering). As a starting point, we focus on engineering α -amylase to obtain enzymes with higher predicted activity. α -Amylase (EC 3.2.1.1) is an enzyme that catalyzes the hydrolysis of the α -d-(1,4)-glucosidic linkages in starch, glycogen, and various oligosaccharides, releasing anomeric products (Fujimoto et al., 1998). Additionally, we study the changes induced in the internal representations of the model comparing the checkpoints of ZymCTRL, a conditional pLM, before and after alignment through direct preference optimization (DPO), to understand the position-dependent patterns learned during RL campaigns.

The contributions of this work are threefold:

- We trained a suite of Sparse Autoencoders on ~ 1 billion tokens from the BRENDA enzyme database (Chang et al., 2020). These SAEs can be applied to diverse downstream tasks, such as explainability or enzyme design.
- We developed a protein engineering workflow by fine-tuning these SAEs on α -amylase deep mutational scanning (DMS) data, identifying features that correlate with fitness through sparse logistic regression. We implemented causal interventions (feature clamping and ablation) with the goal of steering the model toward the desired fitness.
- We analyze how protein language models evolve under RL alignment by applying model diffing, revealing both localized amino acid preference shifts and broader changes in sequence exploration strategies between pre- and post-alignment checkpoints.

2. Methods

2.1. Activity prediction Oracle and Dataset

Following (Schmirler et al., 2024), we trained an activity prediction oracle by fine-tuning ESM-1v (Meier et al., 2021), with LoRA adapters. The model was trained to predict the activity of α -amylase variants using publicly available datasets from the Protein Engineering Tournament GitHub repository (Armer et al., 2023a). Specifically, to predict SAPI values, which represent the ratio of the specific activity of a variant to that of the reference enzyme. Prior to training, we filtered out entries with no recorded activity or with expression below 0.5. The models were trained for 57 epochs using an 80/20 split for training and validation. A batch size of 4 and a learning rate of 3×10^{-4} were applied during training. Learning curves and Spearman correlations are illustrated in Fig. A10.

2.2. SAE architecture and Datasets

We trained a suite of sparse autoencoders on approximately 1 billion tokens from the BRENDA enzyme database (Schomburg et al., 2000), injecting them into the residual stream of ZymCTRL before the attention module. Following best practices, we used the BatchTopK activation function during training, which retains only the top- $k \times b$ activations per batch, where b is the batch size (Fig. 1a).

After pretraining, we fine-tuned each SAE on our Deep Mutational Scanning dataset with a reduced learning rate to prevent overfitting. During training, the batch size was set to 4096, with a learning rate of 3×10^{-4} , using the Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.99$, and an expansion factor of 12. The residual stream dimension is 1280 yielding $1280 \times 12 = 15360$ latents (decoder rows). Layer 26 was chosen based on preliminary results indicating superior performance compared to other insertion points 8.

2.3. Feature Selection and Causal Interventions

To identify the most predictive latent features for enzyme activity, we developed a systematic approach that combined feature selection with causal interventions during sequence generation (Fig. 1b).

Feature Selection Process We pooled position-wise SAE activations into sequence-level vectors by averaging across all sequence positions, resulting in a single feature vector per protein variant. These aggregated vectors were then used to train a Sparse Logistic Regression model using the Sklearn implementation (Pedregosa et al., 2011), with enzyme activity labels (high vs. low activity based on SAPI thresholds) τ as the target variable.

The sparse regularization (L1 penalty) encourages the model to select only the most informative features while setting irrelevant coefficients to zero. The resulting coefficient vector β has as many entries as SAE decoder columns (15,360 latents), with most coefficients being zero due to the sparsity constraint. Features with nonzero coefficients $\beta_i \neq 0$ were identified as predictive of enzyme activity and subsequently used for downstream interventions.

Causal Intervention Strategies We implemented two complementary intervention strategies that modify SAE activations during the forward pass at inference time:

Clamping: For features identified as positively correlated with enzyme activity ($\beta_i > 0$), we set their activations to the maximum observed values in the training distribution whenever these features naturally activated during generation. Formally, if feature i activates with value $f_i > 0$ at position t , we replace it with $f_i^{\text{clamp}} = \max(\mathcal{D}_{\text{train}}^{(i)})$, where $\mathcal{D}_{\text{train}}^{(i)}$ represents all observed activations for feature i in the

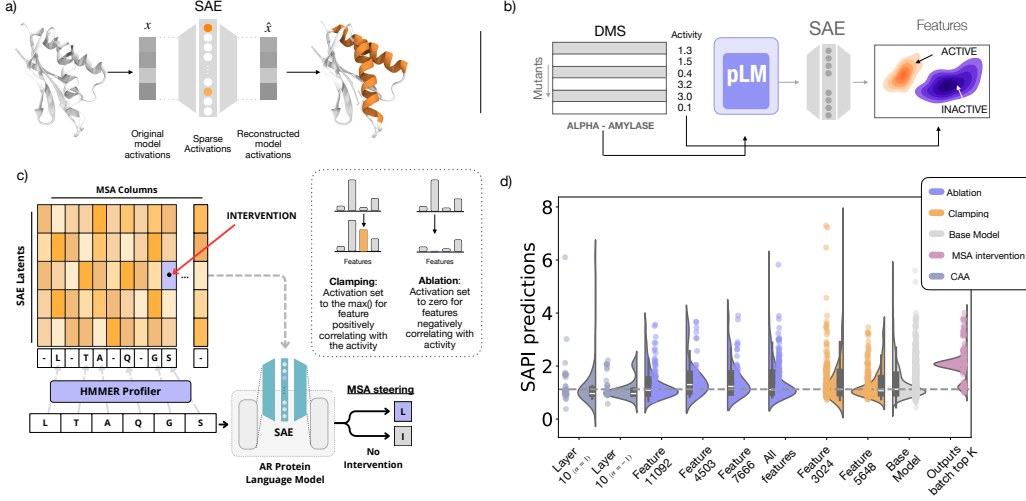


Figure 1. a) Schematic representation of the training process for SAEs. The SAE is inserted between the model’s layers. Embeddings x are passed through the encoder-decoder and reconstructed as \hat{x} , with sparsity enforced in a higher-dimensional space than the input vector. This may provides a more interpretable representation, as learned features can potentially be correlated with observed features. b) Specific application of SAEs for protein engineering, as exemplified in this work. ZymCTRL (pLM) is fed with DMS data, and correlations between learned features and activity measurements are used to interpret and extract relevant features that are then used to steer the model. c) Diagram of the MSA Steering process during inference. To ensure sufficient sequence context for the HMM profiler, interventions begin after approximately 30 amino acids. At each generation step t , the partial sequence is aligned using the profiler to determine the corresponding MSA column. A position-specific intervention is then selected from a pre-computed lookup table designed to increase predicted enzymatic activity, and applied to the SAE latents with reconstruction error preservation. d) Steering is performed through clamping and ablation. The resulting effects reveal an increase in the average predicted activity compared to the base model.

training set (Fig. 1c).

Ablation: For features negatively correlated with enzyme activity ($\beta_i < 0$), we performed ablation by setting their activations to zero whenever they naturally activated. This intervention removes the potentially detrimental influence of these features on sequence generation (Fig. 2c)).

Both intervention types preserve the sparse activation pattern of the SAE while steering the model toward sequences with higher predicted enzymatic activity.

Contrastive Activation Addition (CAA): CAA serves as our baseline steering method (Panickssery et al., 2023). We extract last-token activations $\mathbf{h}_{\text{last}} \in \mathbb{R}^{d_{\text{model}}}$ at layer 25 for all sequences and partition them based on SAPI values relative to threshold τ :

$$\mathcal{H}_{\text{high}} = \{\mathbf{h}_i : \text{SAPI}_i > \tau\} \quad (5)$$

$$\mathcal{H}_{\text{low}} = \{\mathbf{h}_i : \text{SAPI}_i \leq \tau\} \quad (6)$$

The steering vector is computed as the difference between group means:

$$\mathbf{v}_{\text{steer}} = \mu_{\text{high}} - \mu_{\text{low}} = \frac{1}{|\mathcal{H}_{\text{high}}|} \sum_{\mathbf{h} \in \mathcal{H}_{\text{high}}} \mathbf{h} - \frac{1}{|\mathcal{H}_{\text{low}}|} \sum_{\mathbf{h} \in \mathcal{H}_{\text{low}}} \mathbf{h} \quad (7)$$

During generation, this vector is added to layer 25 activations: $\mathbf{h}_{\text{steered}} = \mathbf{h}_{\text{original}} + \alpha \cdot \mathbf{v}_{\text{steer}}$, where α controls intervention strength (Figure A9).

2.4. MSA Steering

In addition to the aforementioned causal intervention methods (ablation, clamping, and CAA steering) that modify the forward pass at inference time, we introduce a novel steering technique: **MSA Steering**. Unlike previous methods that treat all sequence positions equally, MSA Steering designs targeted interventions for each position by leveraging information from previously sampled variants.

Multiple Sequence Alignment (MSA) enables the position-wise aggregation of sequences with different lengths into a unified coordinate system, where each MSA column represents a homologous position across the enzyme variant family.

MSA Steering Methodology The MSA Steering process consists of two phases: preparation and inference-time intervention.

Preparation Phase:

1. Sample and score a library of enzyme variants using the protein language model

2. Extract active SAE latents for each sequence (similar to ablation and clamping)
3. Align the variant library using MSA software (Katoh et al., 2002)
4. Re-index the extracted latents using MSA coordinates into a tensor of shape $n_variants \times msa_columns \times n_latents$
5. Apply Sparse Logistic Regression to identify latent features associated with higher enzymatic activity at each MSA column
6. Train an HMM profiler (Larralde & Zeller, 2023) on the enzyme library to enable alignment of partial sequences during generation

Inference-Time Intervention: During autoregressive sequence generation, interventions are applied after an initial buffer period (30 amino acids) to ensure sufficient sequence context for reliable HMM alignment. The following steps are applied at each time step $t > t_{buffer}$ (where $t_{buffer} \approx 30$ amino acids):

1. **Sample the Next Amino Acid:** At time step t , an amino acid is sampled and added to the growing sequence (s_0, \dots, s_t)
2. **Align with HMM Profiler:** The partial sequence (s_0, \dots, s_t) is aligned using the HMM profiler, mapping it to a homologous region and producing an aligned sequence
3. **Identify MSA Column:** From the aligned sequence, we determine which MSA column corresponds to the current position t
4. **Apply Position-Specific SAE Intervention:** Based on the identified MSA column, we apply the corresponding SAE intervention (e.g., clamping latent features associated with higher enzymatic activity at that specific position)

This approach combines the advantages of model steering with evolutionary information from MSA, enabling more targeted and effective interventions for enzyme design compared to position-agnostic methods.

2.5. Fine tuning and DPO-alignment

ZymCTRL was fine-tuned on 10,398 protein sequences, as detailed in the model card available on Hugging Face (AI4PD/ZymCTRL). Fine-tuning was performed over 28 epochs with a learning rate of 8×10^{-5} on experimental catalytic activity of a DMS library (Armer et al., 2023b). Following fine-tuning, the model was aligned using the

Weighted DPO framework, as described in (Stocco et al., 2024), associating this time to each sequence a measured phenotype. The reward function was defined as the mean of three components: (i) predicted activity, (ii) pLDDT (score, and (iii) TM-score (van Kempen et al., 2023) of the esm-fold (Lin et al., 2023b) predicted protein structure. To mitigate reward hacking and sequence length bias, the final reward was weighted using a Gaussian length penalty centered at 425 residues, the typical length of sequences in the DMS dataset.

2.6. Model Diffing

The pipeline described above maps two global properties of an enzyme variant: its predicted activity and the position-wise pooled SAE activations.

To investigate position-dependent sequence–activity relationships, we compare the next-token probability distributions produced by two checkpoints of our model: the base model and the DPO-aligned model at iteration 3, as it showed the highest reward (Fig. A2)

At each sequence position, we compute the Kullback–Leibler (KL) divergence between the two models’ next-token distributions using the raw, ungapped sequences. For comparison between the two models, we aggregate the KL divergences by aligning the per-position KL divergence by re-indexing based on a multiple sequence alignment (MSA) (Fig. A2), allowing to compare sequences of different lengths. In particular, MSAs of all generated variants are performed using MAFFT (Katoh et al., 2002) with default settings. We then re-index the per-sequence KL divergence scores onto the MSA coordinate frame, so that each divergence value corresponds to a consistent alignment position across variants.

Finally, we select the top MSA columns by average KL divergence. These top-KL positions highlight the residues where the base and DPO-aligned models differ most strongly in their predictive distributions.

3. Experiments and Results

3.1. Steering Interventions for Enzyme Generation

Following Chalnev et al. (Chalnev et al., 2024), we assessed two causal interventions on feature activations during autoregressive generation of enzyme variants. In the *ablation* intervention, whenever a targeted feature naturally activated, its value was set to zero; in the *clamping* intervention, any activation was set to its maximum observed value in the training distribution (Fig. 1c). Both methods relied on reconstruction-error terms from a Sparse Autoencoder to preserve sequence quality. As a baseline, we implemented Contrastive Activation Addition (CAA) (Pan-

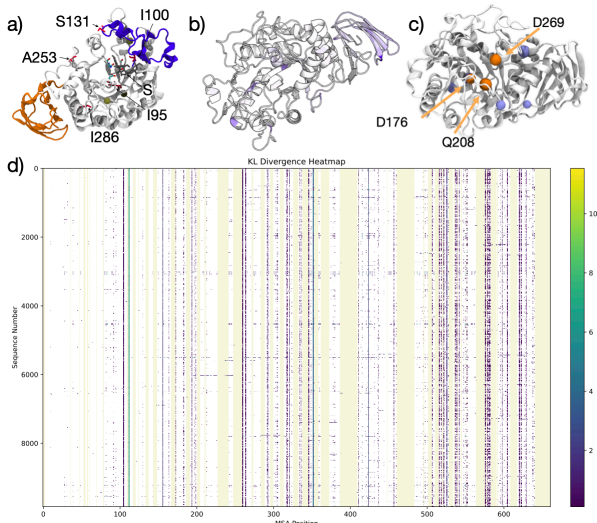


Figure 2. a) Crystal structure (PDB code 1BAG) showing the three domains of α -amylase: Domain C (orange) contains a β -structure with Greek-key motif important for stability, Domain B (purple) captures metal ions and covers the catalytic pocket, and Domain A (white) houses the TIM barrel catalytic domain. The three catalytic residues (two aspartates and one glutamate, with the latter mutated to Gln for crystallization) are depicted, along with the substrate (S) G5 (maltopentaose) in the active site and metal ions (metallic spheres) (Fujimoto et al., 1998). Depicted in magenta the five residues shown in Table 2. b) KL divergence between the DPO and FT model superimposed on the 3d structure of the α -amylase (PDB: 1BAG). c) Position of the transitions described in Table 1 (in purple) respect the catalytic triad (in orange) d) Visualization of the KL divergence for the enzyme library between the FT & DPO models, reindex with MSA column coordinates. The entries in beige correspond to MSA gaps, while white entries correspond to small (≤ 0.1) KL values.

ickssery et al., 2023), which adds a “steering vector” during generation equal to the difference between mean activations of high-activity (> 2.5) versus low-activity (< 2.5) α -amylase classes.

We generated large ensembles under each steering scheme and from the unmodified base model, then predicted their enzymatic activities using our trained oracle. Distributions were compared to the base using the Mann–Whitney U test; only statistically significant shifts were retained for further analysis (Table 1).

MSA Steering produced the largest shift (median +0.94), followed by clamping (+0.13), ablation (+0.07) and CAA (+0.05), confirming that targeted feature manipulations can guide predicted enzyme activity in some cases. Out of the 46 steering interventions tested (17 ablation, 15 clamping, and 13 CAA, 1 MSA Steering, only 12 interventions deviated from the base distribution in a statistically significant

Intervention	Median Predicted Activity	p-value vs. Base
Base (no steering)	1.045	—
Ablation	1.051	0.003
Clamping	1.139	< 0.001
CAA	1.058	0.015
MSA Steering	1.995	< 0.001

Table 1. Median predicted activities and significance of steering interventions compared to the base model.

way. Among these, MSA Steering achieved the most substantial shift, highlighting its potential as a powerful tool for optimizing enzyme activity through model steering.

3.2. Diffing Dynamics During RL Alignment and Interpreting Model Evolution

We applied DPO for three iterations, consisting of less than 0.1% of the compute used in initial pre-training stage (Ferrer & Höcker, 2022)—to align the model towards higher activity. We generated sequences from both the base and DPO-aligned models, performed MSA to re-index per-token similarity metrics, and computed the KL divergence at each MSA position.

Inspection of the highest-divergence positions (Fig. 2) revealed two distinct patterns: sparse, discrete substitutions at key residues (vertical columns on key position that span all the enzyme variants), and broader distributional shifts across contiguous regions of the protein.

3.3. Testing and Quantifying AA Transition Patterns

By exploiting the first type of pattern (discrete substitutions at key residues), we can identify positions whose distribution changed the most through the alignment process with the fitness oracle. From this analysis, five positions (94, 99, 130, 277, 285) exhibited the highest divergence. For each site, we constructed two variant sets: one replacing the wild-type residue with the amino acid favored by the base model, and the other using the DPO-aligned model’s top prediction. All other residues remained unchanged. We then predicted activities for both sets and computed the mean activity difference for each single-point substitution (Table 2).

The I→L substitution at position 285 drove the largest gain (mean +0.946), with I→F at position 99 yielding +0.716. A moderate improvement was observed for S→H at 130 (+0.103), whereas transitions at 94 and 277 were effectively neutral (each +0.010). These results demonstrate that a handful of targeted amino acid changes can recapitulate most of the alignment-induced activity enhancements.

Residue Position	AA Transition	Δ Mean Activity
94	I \rightarrow L	0.010
99	I \rightarrow F	0.716
130	S \rightarrow H	0.103
277	A \rightarrow L	0.010
285	I \rightarrow L	0.946

Table 2. Activity shifts for single-point mutations informed by base vs. & DPO model preferences. The residue positions are aligned with the WT Enzyme Variant 1BAG

4. Discussion and Limitations

Reverse-engineering to make neural networks human-interpretable is the aim of mechanistic interpretability (Olah et al., 2018; Meng et al., 2022; Nanda et al., 2023). A key challenge of mechanistic interpretability is identifying the correct units of analysis, that are ideally canonical (irreducible, indivisible, and complete) (Leask et al., 2025). Due to their properties, SAEs offer intriguing possibilities for interpretability research.

In this work, we explored the application of SAEs in the context of a protein engineering campaign. Specifically, we trained SAEs and extracted features that correlate with an external oracle trained to predict enzyme activity. By ablating and clamping targetted activations, we observed it is possible to deviate the base model distribution, although the effect of a single intervention at a time remains modest. We introduce a new intervention strategy grounded in the multiple sequence alignment (MSA) of the engineered protein, projecting SAE latents into MSA space to enable precise and effective interventions that leverage evolutionarily conserved features. We also computed KL divergences between base and aligned models, to investigate how RL alters the model’s internal representations. Through this process, we were able to capture fine-grained differences and identify how individual mutations contributed to measurable improvements in generated sequences.

Interestingly, most of the substitutions occur between chemically similar residues (I \rightarrow L), whereas the serine to histidine transition (S \rightarrow H) represent a shift from a polar, uncharged side chain to a positively charged residue. All of these mutations lie in the enzyme third and forth coordination shells, highlighting once again, that the impact of mutations and the underlying allosteric network effects are extremely not obvious to decipher (Osuna, 2021; Gu et al., 2023). In this context, molecular dynamics simulations can provide valuable insights in the mechanistic effects of these transitions, which have been linked to improved predicted activity. Interestingly, these transitions are also present in the DMS data, where they strongly correlate with increased activity. This finding confirms that the model is properly aligned and suggests the potential for using this tool to extrapolate key

elements solely from model-based investigations.

In future work, we envision (1) further research on specific circuits, at different MSA-positions during sequential DPO rounds (2) testing base and steered designs experimentally to validate our approach.

Acknowledgements

This work is supported by ERC grant ATHENA, Grant Agreement 101165231.

Contributions

G.B.C conceived the work, trained the models, analyzed the data, and wrote the manuscript. F.S trained the aligned models and wrote the manuscript. N.F analyzed the data and wrote the manuscript. The three authors discussed the results and supervised the work.

References

- Adams, E., Bai, L., Lee, M., Yu, Y., and AlQuraishi, M. From mechanistic interpretability to mechanistic biology: Training, evaluating, and interpreting sparse autoencoders on protein language models. *bioRxiv*, pp. 2025–02, 2025.
- Armer, C., Kane, H., Cortade, D., Estell, D., Yusuf, A., Sanka, R., Redestig, H., Brunette, T., Kelly, P., and DeBenedictis, E. The Protein Engineering Tournament: An Open Science Benchmark for Protein Modeling and Design, 2023a. URL <https://arxiv.org/abs/2309.09955>. Version Number: 2.
- Armer, C., Kane, H., Cortade, D., Estell, D., Yusuf, A., Sanka, R., Redestig, H., Brunette, T., Kelly, P., and DeBenedictis, E. The protein engineering tournament: An open science benchmark for protein modeling and design, 2023b. URL <https://arxiv.org/abs/2309.09955>.
- Bhatnagar, A., Jain, S., Beazer, J., Curran, S. C., Hoffnagle, A. M., Ching, K., Martyn, M., Nayfach, S., Ruffolo, J. A., and Madani, A. Scaling unlocks broader generation and deeper functional understanding of proteins, April 2025. URL <http://biorxiv.org/lookup/doi/10.1101/2025.04.15.649055>.
- Bussmann, B., Leask, P., and Nanda, N. Batchtopk sparse autoencoders. *arXiv preprint arXiv:2412.06410*, 2024.
- Chalnev, S., Siu, M., and Conmy, A. Improving steering vectors by targeting sparse autoencoder features. *arXiv preprint arXiv:2411.02193*, 2024.
- Chang, A., Jeske, L., Ulbrich, S., Hofmann, J., Koblit, J., Schomburg, I., Neumann-Schaal, M., Jahn, D., and

- Schomburg, D. Brenda, the elixir core data resource in 2021: new developments and updates. *Nucleic Acids Research*, 49(D1):D498–D508, 11 2020. ISSN 0305-1048. doi: 10.1093/nar/gkaa1025. URL <https://doi.org/10.1093/nar/gkaa1025>.
- Ferruz, N. and Höcker, B. Controllable protein design with language models. *Nature Machine Intelligence*, 4(6): 521–532, 2022.
- Fujimoto, Z., Takase, K., Doui, N., Momma, M., Matsumoto, T., and Mizuno, H. Crystal structure of a catalytic-site mutant α -amylase from *Bacillus subtilis* complexed with maltopentaose. *Journal of Molecular Biology*, 277(2):393–407, March 1998. ISSN 00222836. doi: 10.1006/jmbi.1997.1599. URL <https://linkinghub.elsevier.com/retrieve/pii/S0022283697915990>.
- Garcia, E. N. V. and Ansuini, A. Interpreting and Steering Protein Language Models through Sparse Autoencoders, 2025. URL <https://arxiv.org/abs/2502.09135>. Version Number: 1.
- Gu, J., Xu, Y., and Nie, Y. Role of distal sites in enzyme engineering. *Biotechnology Advances*, 63:108094, March 2023. ISSN 07349750. doi: 10.1016/j.biotechadv.2023.108094. URL <https://linkinghub.elsevier.com/retrieve/pii/S0734975023000010>.
- Johnson, S. R., Fu, X., Viknander, S., Goldin, C., Monaco, S., Zelezniak, A., and Yang, K. K. Computational Scoring and Experimental Evaluation of Enzymes Generated by Neural Networks, March 2023. URL <http://biorxiv.org/lookup/doi/10.1101/2023.03.04.531015>.
- Katoh, K., Misawa, K., Kuma, K.-i., and Miyata, T. Mafft: a novel method for rapid multiple sequence alignment based on fast fourier transform. *Nucleic acids research*, 30(14):3059–3066, 2002.
- Larralde, M. and Zeller, G. Pyhmmer: a python library binding to hmmer for efficient sequence analysis. *Bioinformatics*, 39(5):btad214, 04 2023. ISSN 1367-4811. doi: 10.1093/bioinformatics/btad214. URL <https://doi.org/10.1093/bioinformatics/btad214>.
- Leask, P., Bussmann, B., Pearce, M., Bloom, J., Tigges, C., Moubayed, N. A., Sharkey, L., and Nanda, N. Sparse autoencoders do not find canonical units of analysis, 2025. URL <https://arxiv.org/abs/2502.04878>.
- Lin, Z., Akin, H., Rao, R., Hie, B., Zhu, Z., Lu, W., Smetanin, N., Verkuil, R., Kabeli, O., Shmueli, Y., Dos Santos Costa, A., Fazel-Zarandi, M., Sercu, T., Candido, S., and Rives, A. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, March 2023a. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.ade2574. URL <https://www.science.org/doi/10.1126/science.ade2574>.
- Lin, Z., Akin, H., Rao, R., Hie, B., Zhu, Z., Lu, W., Smetanin, N., Verkuil, R., Kabeli, O., Shmueli, Y., et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637): 1123–1130, 2023b.
- Madani, A., Krause, B., Greene, E. R., Subramanian, S., Mohr, B. P., Holton, J. M., Olmos, J. L., Xiong, C., Sun, Z. Z., Socher, R., Fraser, J. S., and Naik, N. Large language models generate functional protein sequences across diverse families. *Nature Biotechnology*, 41(8):1099–1106, August 2023. ISSN 1087-0156, 1546-1696. doi: 10.1038/s41587-022-01618-2. URL <https://www.nature.com/articles/s41587-022-01618-2>.
- Marks, S., Rager, C., Michaud, E. J., Belinkov, Y., Bau, D., and Mueller, A. Sparse feature circuits: Discovering and editing interpretable causal graphs in language models. *arXiv preprint arXiv:2403.19647*, 2024.
- Meier, J., Rao, R., Verkuil, R., Liu, J., Sercu, T., and Rives, A. Language models enable zero-shot prediction of the effects of mutations on protein function. *bioRxiv*, July 2021. doi: 10.1101/2021.07.09.450648. URL <http://dx.doi.org/10.1101/2021.07.09.450648>.
- Meng, K., Bau, D., Andonian, A., and Belinkov, Y. Locating and editing factual associations in gpt. *Advances in neural information processing systems*, 35:17359–17372, 2022.
- Munsamy, G., Lindner, S., Lorenz, P., and Ferruz, N. Zymctrl: a conditional language model for the controllable generation of artificial enzymes. In *NeurIPS machine learning in structural biology workshop*, 2022.
- Nanda, N., Rajamanoharan, S., Kramar, J., and Shah, R. Fact finding: Attempting to reverse-engineer factual recall on the neuron level. In *Alignment Forum*, pp. 6, 2023.
- Olah, C., Satyanarayan, A., Johnson, I., Carter, S., Schubert, L., Ye, K., and Mordvintsev, A. The building blocks of interpretability. *Distill*, 3(3):e10, 2018.
- Osuna, S. The challenge of predicting distal active site mutations in computational enzyme design. *WIREs Comput Mol Sci*, 11(3):e1502, May 2021. ISSN 1759-0876, 1759-0884. doi: 10.1002/wcms.1502. URL <https://wires.onlinelibrary.wiley.com/doi/10.1002/wcms.1502>.

- Panickssery, N., Gabrieli, N., Schulz, J., Tong, M., Hubinger, E., and Turner, A. M. Steering llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*, 2023.
- Parsan, N., Yang, D. J., and Yang, J. J. Towards interpretable protein structure prediction with sparse autoencoders. *arXiv preprint arXiv:2503.08764*, 2025.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Schmirler, R., Heinzinger, M., and Rost, B. Fine-tuning protein language models boosts predictions across diverse tasks. *Nature Communications*, 15(1):7407, 2024.
- Schomburg, I., Hofmann, O., Baensch, C., Chang, A., and Schomburg, D. Enzyme data and metabolic information: Brenda, a resource for research in biology, biochemistry, and medicine. *Gene Function & Disease*, 1(3-4):109–118, 2000.
- Simon, E. and Zou, J. Interplm: Discovering interpretable features in protein language models via sparse autoencoders. *bioRxiv*, pp. 2024–11, 2024.
- Stocco, F., Artigues-Lleixa, M., Hunklinger, A., Widatalla, T., Guell, M., and Ferruz, N. Guiding generative protein language models with reinforcement learning. *arXiv preprint arXiv:2412.12979*, 2024.
- van Kempen, M., Kim, S. S., Tumescheit, C., Mirdita, M., Lee, J., Gilchrist, C. L. M., Söding, J., and Steinegger, M. Fast and accurate protein structure search with foldseek. *Nature Biotechnology*, 42(2): 243–246, May 2023. ISSN 1546-1696. doi: 10.1038/s41587-023-01773-0. URL <http://dx.doi.org/10.1038/s41587-023-01773-0>.
- Yang, J., Bhatnagar, A., Ruffolo, J. A., and Madani, A. Conditional Enzyme Generation Using Protein Language Models with Adapters, 2024. URL <https://arxiv.org/abs/2410.03634>. Version Number: 1.

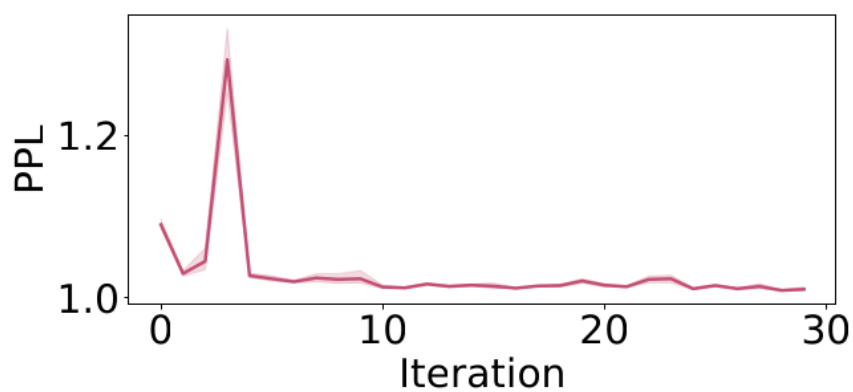


Figure 3. Average sequence perplexity across sequential DPO rounds.

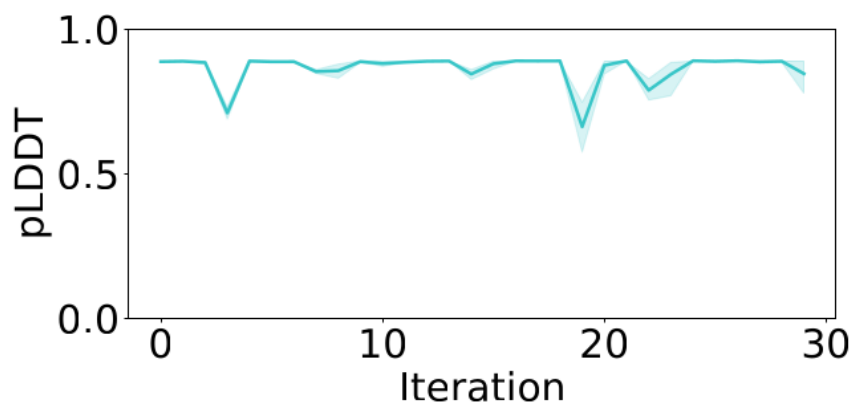


Figure 4. Average pLDDT, as measured by ESMFold, across sequential DPO rounds.

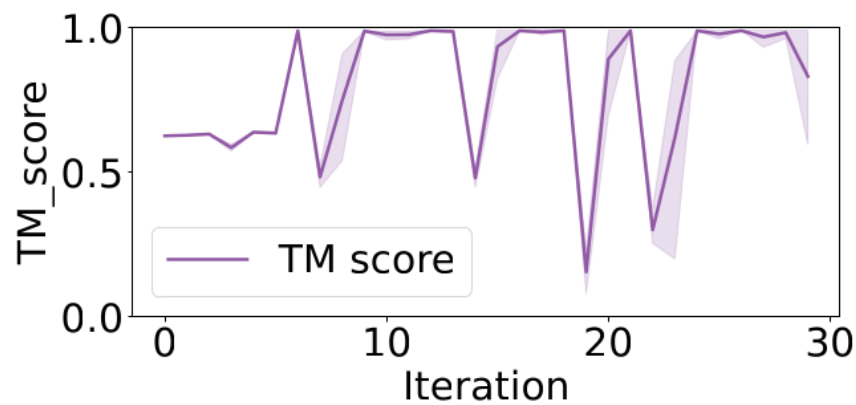


Figure 5. Average TM-score between the reference enzyme and ESMFold-predicted structures during sequential DPO rounds.

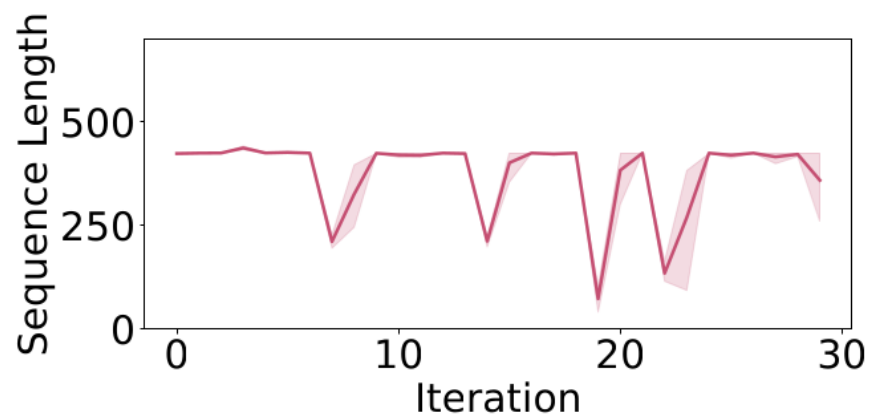


Figure 6. Average sequence length across sequential DPO rounds.

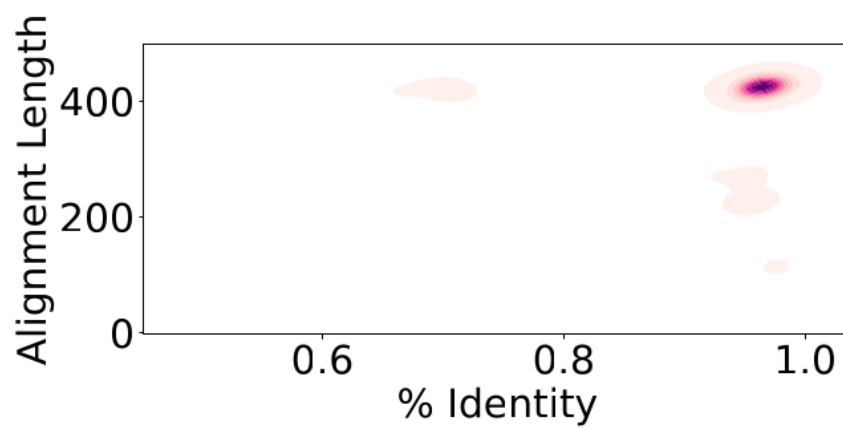


Figure 7. Global distribution of sequence length versus alignment length.

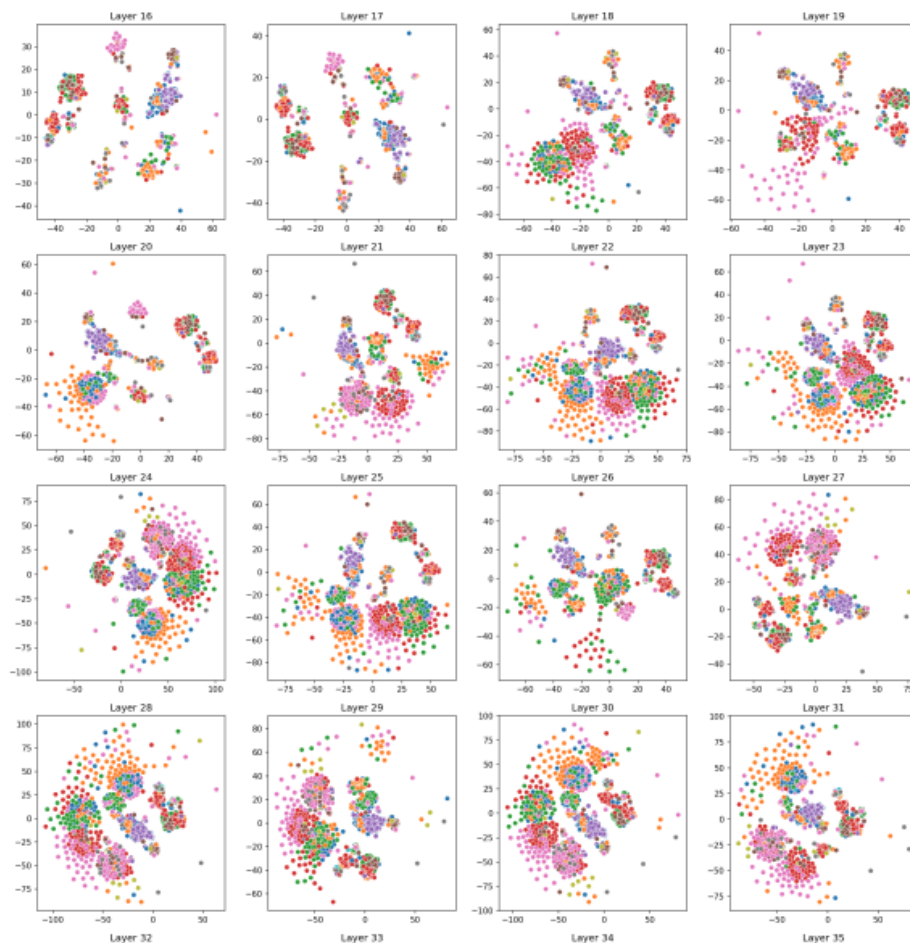


Figure 8. TSNE visualization of the embeddings of DPO sequences at different layers

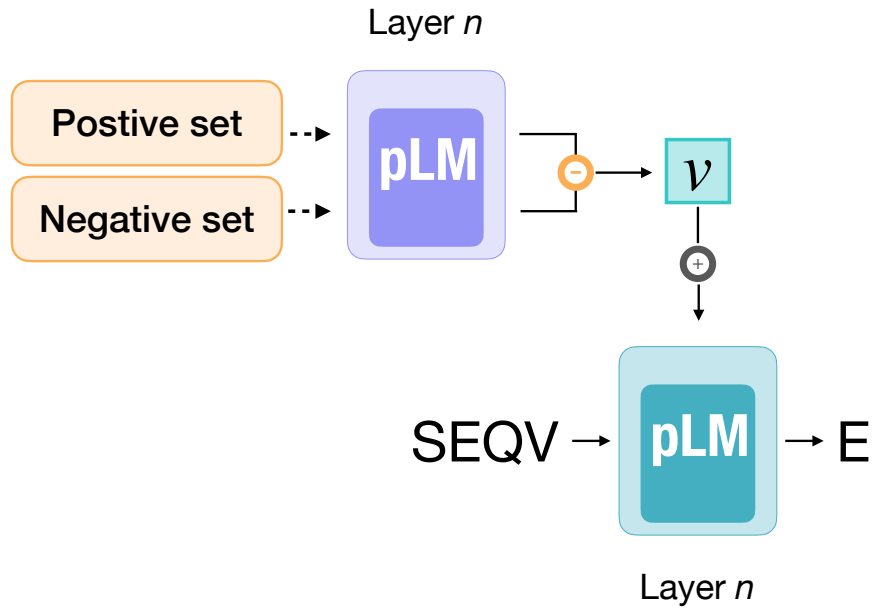


Figure 9. Overview of Contrastive Activation Addition (CAA) steering. The steering vector is computed as the mean difference in activations between positive (desired property) and negative (undesired property) protein sets. During generation, the vector v is added to the model activations at the same layer, the residual stream results thus modified by adding the vector v times an scalar α , that controls the strength of steering.

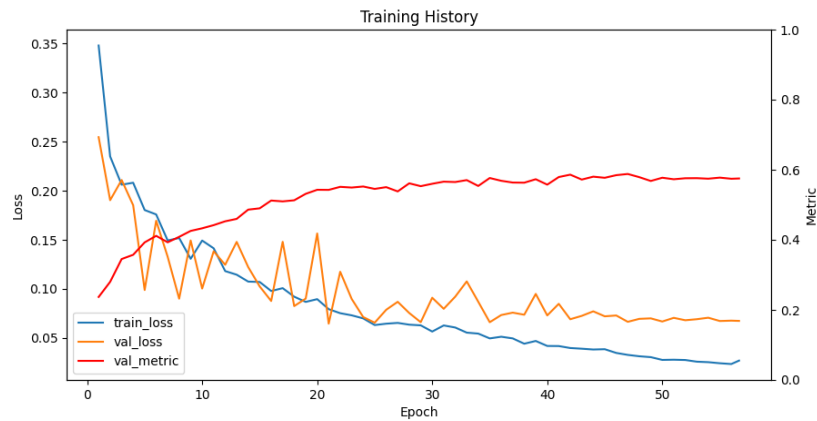


Figure 10. Training curve of esm-1v with Lora Adapter, as reported in Chalnev et al. (Chalnev et al., 2024)