

ProtSAE: Disentangling and Interpreting Protein Language Models via Semantically-Guided Sparse Autoencoders

Xiangyu Liu¹, Haodi Lei¹, Yi Liu¹, Yang Liu¹, Wei Hu^{1,2,*}

¹ State Key Laboratory for Novel Software Technology, Nanjing University, China

² National Institute of Healthcare Data Science, Nanjing University, China

{xyl.nju, haodilei, yiliu07.nju, yliu20.nju}.nju@gmail.com, whu@nju.edu.cn

Abstract

Sparse Autoencoder (SAE) has emerged as a powerful tool for mechanistic interpretability of large language models. Recent works apply SAE to protein language models (PLMs), aiming to extract and analyze biologically meaningful features from their latent spaces. However, SAE suffers from semantic entanglement, where individual neurons often mix multiple nonlinear concepts, making it difficult to reliably interpret or manipulate model behaviors. In this paper, we propose a semantically-guided SAE, called ProtSAE. Unlike existing SAE which requires annotation datasets to filter and interpret activations, we guide semantic disentanglement during training using both annotation datasets and domain knowledge to mitigate the effects of entangled attributes. We design interpretability experiments showing that ProtSAE learns more biologically relevant and interpretable hidden features compared to previous methods. Performance analyses further demonstrate that ProtSAE maintains high reconstruction fidelity while achieving better results in interpretable probing. We also show the potential of ProtSAE in steering PLMs for downstream generation tasks.

1 Introduction

In recent years, protein language models (PLMs) (Lin et al. 2023a; Nijkamp et al. 2023) have developed rapidly and been widely applied to downstream tasks including protein function prediction (Lin et al. 2024), structural modeling (Lin et al. 2023b), and protein design (Ferruz and Höcker 2022). However, the internal mechanisms of PLMs remain largely unknown (Garcia and Ansuini 2025).

For protein engineering, it is important to understand how latent features map to biological concepts, such as binding pockets, post-translational modifications, or fold families. Such analysis facilitates the identification of spurious correlations and biases, enhancing both performance and robustness of PLMs. Moreover, it allows for the extraction of latent relationships among protein characteristics, providing meaningful insights that can inform and support biological research (Zhang et al. 2024). Early works have attempted to analyze protein models, exploring the relationships between attention mechanisms and amino acids (Vig et al. 2021), as well as identifying neurons associated with certain biological concepts (Nori, Singireddy, and Have 2023a).

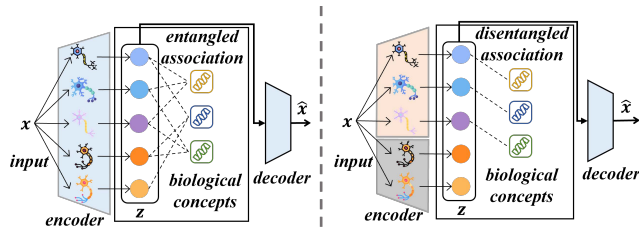


Figure 1: Illustration of SAE semantic entanglement (left): individual neurons conflate multiple biological concepts, and semantic disentanglement (right): each defined neuron maps to a single biological concept.

Recent studies (Simon and Zou 2024; Garcia and Ansuini 2025; Adams et al. 2025) have begun applying sparse autoencoder (SAE) to PLMs and observed the emergence of features associated with various biological concepts. SAE is an effective tool for understanding and explaining the internal representations of PLMs. Based on the assumption of linear feature superposition (Yun et al. 2021), it decomposes the hidden representations to extract sparse features. These features can be further analyzed for their correlations with specific concepts, enabling interpretability. Furthermore, the sparse features can be selectively activated to steer the generation along the directions of relevant concepts.

A typical SAE requires an annotation dataset after training to interpret the learned features (Simon and Zou 2024). This annotation contains concepts of interest, and post-hoc correlation analysis is often performed to establish the relationship between the features in SAE and these concepts. However, SAE suffers from the problem of semantic entanglement: *individual neurons often conflate multiple concepts* (Joshi et al. 2025). This entanglement results in ambiguous interpretations of the learned features. As illustrated in Figure 1, each neuron is likely to be simultaneously associated with multiple, semantically divergent concepts. Consequently, identifying the true meaning of the given feature and using it to steer the generation becomes challenging, undermining the interpretability of the model.

In this paper, we propose ProtSAE, which incorporates semantic guidance into the SAE training to disentangle semantic features. First, we leverage the semantic annotations used

*Corresponding author

for post-hoc interpretation of SAE features to constrain the relationship between defined activations and specific concepts during training. We also use forced activations and feature rescaling to ensure that the defined activations effectively participate in reconstruction with high fidelity. Second, considering the rich prior knowledge in the protein domain, where concepts are not mutually independent, we incorporate ELEmbeddings (Kulmanov et al. 2019) to model potential logical constraints among concepts, e.g., subsumption and conjunction. The constraints are integrated into the training process of ProtSAE to enhance the interpretability and semantic consistency of the learned features.

We construct interpretability experiments to demonstrate that ProtSAE effectively captures biologically meaningful features in PLMs, such as molecular functions, biological processes, and binding sites. Compared with features annotated from the typical SAE, the defined neurons in ProtSAE learn more accurate, disentangled representations that are more tightly aligned with protein structures. Furthermore, performance analyses show that ProtSAE preserves richer semantic information related to protein concepts. Under varying levels of sparsity, it consistently achieves stronger performance on protein function prediction (Kulmanov et al. 2024; Kulmanov and Hoehndorf 2022) while maintaining high reconstruction fidelity. Finally, through targeted activation steering across various biological concepts, we demonstrate that the semantic features learned by ProtSAE can effectively guide PLM outputs toward desired functional outcomes—validating both the quality of the learned representations and their potential for precise model control.

The main contributions of this paper are listed as follows:

- We propose ProtSAE, a novel semantically-guided SAE that can disentangle complex protein features, yielding features in PLMs more strongly aligned with biological concepts.
- We introduce protein domain knowledge into ProtSAE training to learn the logical constraints among concepts, and apply forced activations and feature rescaling to ensure that defined activations effectively participate in reconstruction with high fidelity.
- We conduct extensive experiments and analyses. Interpretability experiments demonstrate that ProtSAE captures more interpretable features that are closely aligned with biological concepts. Detailed performance analyses show that ProtSAE consistently outperforms baselines across varying levels of sparsity while maintaining high reconstruction fidelity. Steering experiments reveal that ProtSAE enables effective interventions across diverse biological concepts.

2 Related Work

Mechanistic interpretability and SAE. Mechanistic interpretability aims to understand how neural networks produce outputs based on the internal algorithms that they have learned (Olah et al. 2020). Previous works explore the computation subgraphs responsible for specific tasks (Shi et al. 2024; Dunefsky, Chlenski, and Nanda 2024; Wang

et al. 2023), and analyze the behaviors within large language models (LLMs) (Makelov 2024; Miller, Chughtai, and Saunders 2024; Makelov et al. 2024). One prominent line of analysis focuses on identifying and studying sparse linear features within LLMs (Todd et al. 2024). Based on the assumption of linear feature superposition (Elhage et al. 2022), some works decompose language model activations and use them to intervene in the model’s behavior (Yun et al. 2021; Tamkin, Taufeeque, and Goodman 2024). Recent scaling efforts demonstrate the viability of SAE across LLMs, from Claude 3 Sonnet (Paulo et al. 2024) to GPT-4 (Gao et al. 2024), with extensions to multi-modal LLMs as well (Pach et al. 2025). Several works propose architectural improvements to SAE to mitigate feature shrinkage and improve reconstruction fidelity (Wright and Sharkey 2024; Rajamanoharan et al. 2024). Others focus on developing comprehensive evaluation frameworks for SAE (Gallifant et al. 2025) and exploring generating more informative explanations for activated features with additional datasets or LLMs (Wu et al. 2025a,b).

Interpretability in PLMs. Early works in PLMs show that attention maps can capture structural and functional signals, including amino acid interactions (Vig et al. 2021), protein contacts (Rao et al. 2021), and functional sites like binding pockets and allosteric regions (Kannan, Hie, and Kim 2024; Dong et al. 2024). CB-pLM (Ismail et al. 2025) trains PLMs with a concept bottleneck layer for better understanding and controlling PLMs’ generation. Recent studies explore how high-level conceptual knowledge is internally represented in the components of PLMs (Nori, Singireddy, and Have 2023b). SAE is used to decompose latent activations and reveal links between biological concepts and structural features (Simon and Zou 2024; Garcia and Ansuini 2025; Adams et al. 2025). These studies also show that editing related activations can influence or steer protein sequence generation. However, features from SAE often entangle multiple concepts, making interpretation unclear. To address this, we introduce semantic guidance during the SAE training by linking specific activations to biological concepts, leading to more interpretable features.

3 Overview

SAE architectures. SAE is designed to learn sparse representations of high-dimensional inputs by encouraging only a small subset of neurons to be activated. It is widely used to extract interpretable and localized features, particularly in the context of mechanistic interpretability for deep models.

Given an input vector $\mathbf{x} \in \mathbb{R}^d$, the encoder maps it to the latent activations $\mathbf{z} \in \mathbb{R}^n$ using a linear transformation followed by the ReLU activation:

$$\mathbf{z} = \text{ReLU}(\mathbf{W}_{\text{enc}}(\mathbf{x} - \mathbf{b}_{\text{dec}}) + \mathbf{b}_{\text{enc}}), \quad (1)$$

where $\mathbf{W}_{\text{enc}} \in \mathbb{R}^{n \times d}$, $\mathbf{b}_{\text{enc}} \in \mathbb{R}^n$, and $\mathbf{b}_{\text{dec}} \in \mathbb{R}^d$ are learnable parameters. To enable sparsity in \mathbf{z} , n is typically set much larger than d ($n \gg d$). The decoder reconstructs the input via a linear transformation:

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

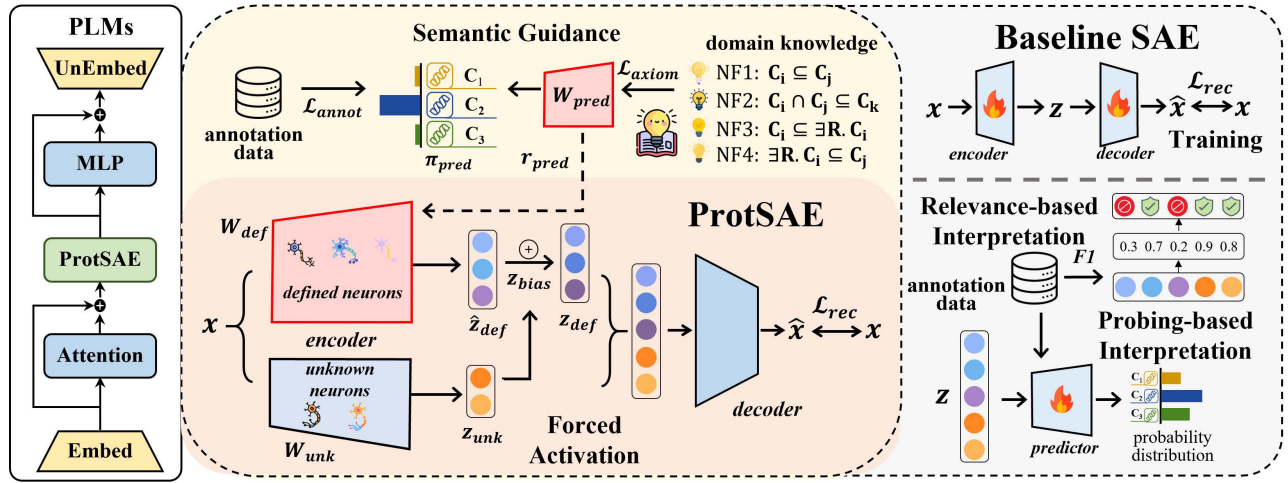


Figure 2: An overview of ProtSAE (left) and the baseline SAE (right). In the baseline SAE, annotation data is used post hoc to interpret learned features via relevance-based and probing-based methods. In contrast, ProtSAE incorporates semantic guidance during training by leveraging annotation data and protein domain knowledge to achieve semantic disentanglement. It also uses forced activations and feature rescaling to learn meaningful features while preserving reconstruction fidelity.

where $\mathbf{W}_{\text{dec}} \in \mathbb{R}^{d \times n}$. The training objective encourages accurate reconstruction and sparsity in \mathbf{z} :

$$\mathcal{L} = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + \lambda \|\mathbf{z}\|_1, \quad (3)$$

where $\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$ is the reconstruction MSE loss, $\|\mathbf{z}\|_1$ imposes an L1 penalty to encourage sparsity, and λ is a tunable hyperparameter to balance the two terms.

Interpreting SAE activations. As shown in Figure 2, we follow prior works and use an annotation dataset that maps proteins to biological concepts to interpret the learned SAE features via two approaches: (1) *Relevance-based interpretation*. Following previous work (Garcia and Ansuini 2025), we compute the activation levels of each feature across annotated proteins. Based on the annotation, we calculate a relevance score (e.g., F1-score) between a feature and a target concept, and use the relevant concept to interpret the feature. Detailed formulations are described in Appendix C.4. (2) *Probing-based interpretation*. We train linear probing classifiers on SAE activations using the annotation dataset (Simon and Zou 2024; Gurnee et al. 2023). It aims to detect the presence of specific concepts within the learned features through supervised training.

4 Method

Conventionally, the encoder of SAE serves two purposes: (1) *Determining which features should be active*. To disentangle semantics in activated features, each defined neuron should be selectively activated only by proteins associated with a specific concept, while remaining inactive for unrelated protein sequences. This constraint helps prevent entangled semantics within the same neuron. Based on this intuition, we introduce semantic guidance from the annotation data and domain knowledge to learn such semantically selective activations during training. (2) *Estimating the magnitude of active features to support faithful reconstruction*. Although

the magnitude of each active feature should be determined by the reconstruction training, we must ensure that the feature directions associated with the predefined concepts effectively contribute to reconstruction. This guarantees that steering the encoder’s activation yields consistent and interpretable effects on the model’s behavior.

4.1 Guiding SAE with Annotation Data

We adopt the TopK-SAE as the backbone, where only the top- K neurons with the highest activations are used in reconstruction. We define the encoder as follows:

$$\mathbf{z} = \text{TopK}(\mathbf{W}_{\text{enc}}(\mathbf{x} - \mathbf{b}_{\text{dec}}) + \mathbf{b}_{\text{enc}}), \quad (4)$$

where $\mathbf{W}_{\text{enc}} \in \mathbb{R}^{n \times d}$, $\mathbf{b}_{\text{enc}} \in \mathbb{R}^n$, and $\mathbf{b}_{\text{dec}} \in \mathbb{R}^d$. $\text{TopK}(\cdot)$ retains only the top K largest activations, zeroing out the rest. Suppose that there are m defined concepts of interest. We aim for m corresponding activations \mathbf{z}_{def} to accurately represent these concepts, while retaining the remaining $n - m$ activations \mathbf{z}_{unk} to capture unknown semantic concepts. We partition the activation \mathbf{z} , weight matrix \mathbf{W}_{enc} , and \mathbf{b}_{enc} into two components:

$$\mathbf{z} = \begin{bmatrix} \mathbf{z}_{\text{def}} \\ \mathbf{z}_{\text{unk}} \end{bmatrix}, \mathbf{W}_{\text{enc}} = \begin{bmatrix} \mathbf{W}_{\text{def}} \\ \mathbf{W}_{\text{unk}} \end{bmatrix}, \mathbf{b}_{\text{enc}} = \begin{bmatrix} \mathbf{b}_{\text{def}} \\ \mathbf{b}_{\text{unk}} \end{bmatrix}, \quad (5)$$

where $\mathbf{W}_{\text{def}} \in \mathbb{R}^{m \times d}$ and $\mathbf{b}_{\text{def}} \in \mathbb{R}^m$ corresponds to m defined activations \mathbf{z}_{def} aligned with predefined concepts, $\mathbf{W}_{\text{unk}} \in \mathbb{R}^{(n-m) \times d}$ and $\mathbf{b}_{\text{unk}} \in \mathbb{R}^{n-m}$ capture the remaining activations \mathbf{z}_{unk} .

Semantic disentanglement. To guide the semantic disentanglement of specific activations, we introduce a concept predictor. It learns to estimate the presence of each predefined concept from the input. Let $\mathbf{W}_{\text{pred}} \in \mathbb{R}^{m \times d}$ denote its weight matrix, the prediction probability of defined activations is computed as follows:

$$\pi_{\text{pred}} = \sigma(\mathbf{W}_{\text{pred}}(\mathbf{x} - \mathbf{b}_{\text{dec}}) + \mathbf{b}_{\text{pred}}) \in (0, 1)^m, \quad (6)$$

where $\sigma(\cdot)$ is the sigmoid function and $\mathbf{b}_{\text{pred}} \in \mathbb{R}^m$ is the prediction bias. We train this predictor on the available annotation data using a binary cross-entropy loss:

$$\mathcal{L}_{\text{annot}} = \text{CrossEntropy}(\pi_{\text{pred}}, y), \quad (7)$$

where $y \in \{0, 1\}^m$ is the binary annotation vector indicating which semantic concepts are present.

We assume that \mathbf{W}_{def} (used for reconstruction) and \mathbf{W}_{pred} (used for prediction) encode the same underlying semantic meanings. Thus, they should share the same projection directions. To achieve this, we treat \mathbf{W}_{def} as a rescaled version of \mathbf{W}_{pred} , and tie their weights as follows:

$$\mathbf{W}_{\text{def}} = \mathbf{W}_{\text{pred}}^{\text{detach}} \cdot \exp(\mathbf{r}_{\text{pred}}), \quad (8)$$

where $\mathbf{r}_{\text{pred}} \in \mathbb{R}^m$ is a learnable scaling vector and \cdot denotes row-wise multiplication, where each row i of $\mathbf{W}_{\text{pred}}^{\text{detach}}$ is scaled by $\exp(\mathbf{r}_{\text{pred}}[i])$. The `detach` indicates that gradients from the reconstruction loss are prevented from updating \mathbf{W}_{pred} . This formulation ensures that \mathbf{W}_{def} retains the semantic directionality learned from supervision, while its magnitude can adapt to improve reconstruction. The exponential guarantees positivity and allows smooth multiplicative modulation.

Forced activation. Using the encoder, we compute the semantic and unsupervised activations as

$$\mathbf{z}_{\text{unk}} = \text{TopK}(\mathbf{W}_{\text{unk}}(\mathbf{x} - \mathbf{b}_{\text{dec}}) + \mathbf{b}_{\text{unk}}), \quad (9)$$

$$\hat{\mathbf{z}}_{\text{def}} = \mathbf{W}_{\text{def}}(\mathbf{x} - \mathbf{b}_{\text{dec}}) + \mathbf{b}_{\text{def}}, \quad (10)$$

where $\hat{\mathbf{z}}_{\text{def}}$ denotes the pre-activation output for the defined concept neurons before any sparsity constraints. In practice, we observe that the reconstruction tends to rely more on the entangled, unsupervised activations \mathbf{z}_{unk} , which diminishes the contribution of concept-specific activations \mathbf{z}_{def} . To mitigate this issue, we introduce a semantic bias that encourages the activations associated with predicted concepts to contribute more strongly to reconstruction:

$$\mathbf{z}_{\text{bias}} = \mathbb{1}_{\pi_{\text{pred}} > 0} \cdot \text{ReLU}(\text{mean}(\mathbf{z}_{\text{unk}}) - \hat{\mathbf{z}}_{\text{def}}), \quad (11)$$

$$\mathbf{z}_{\text{def}} = \hat{\mathbf{z}}_{\text{def}} + \mathbf{z}_{\text{bias}}. \quad (12)$$

Here, $\mathbb{1}_{\pi_{\text{pred}} > 0} \in \{0, 1\}^m$ denotes an indicator function, marking whether each semantic concept is predicted to be present (i.e., $\pi_{\text{pred}} > 0$) and \cdot denotes element-wise multiplication. For such concepts, we enforce the corresponding activation to be no less than the average activation of \mathbf{z}_{unk} , preventing the model from ignoring semantically meaningful features during reconstruction. This bias effectively forces the features aligned with known semantics to participate in encoding, enhancing both interpretability and task relevance.

4.2 Guiding SAE with Domain Knowledge

In protein sequence modeling, biological concepts are often semantically interdependent. The relationship of some concepts can be defined with logical constraints including “is-a”, “part-of”, “regulates”, and other relations. These axioms establish a stable, expert-curated structure over biological knowledge, dictating how concepts relate and compose.

Therefore, aligning latent directions in SAE’s hidden space with concept semantic relationships can enhance the detection and disentanglement of meaningful biological concepts.

To achieve this, we incorporate domain knowledge into SAE using ELEmbeddings (Kulmanov et al. 2019). It is an ontology representation learning method, which represents each concept as a hypersphere in the embedding space, and encodes logical axioms as constraints on the positions and relationships between these regions. Given a concept c_i , the prediction probability of a protein p can be modeled using ELEmbeddings as

$$y'_i = \sigma \left(f_\eta(p)^\top \cdot (f_\eta(hF) + f_\eta(c_i)) + r_\eta(c_i) \right), \quad (13)$$

where $f_\eta(\cdot)$ is the projection function into the semantic embedding space, hF is the hasFunction relation, and $r_\eta(c_i) \in \mathbb{R}_{>0}$ is a learned radius bias. In Appendix A.1, we prove that Eq. (6) is structurally equivalent to Eq. (13). Thus, from the perspective of ontology representation learning, the weight matrix \mathbf{W}_{pred} learned on the LLM latent space can be interpreted as an ontology embedding with relational biases.

We adopt four normalized axiom forms supported by ELEmbeddings, each corresponding to a specific type of logical relation commonly found in ontologies:

- **NF1.** Subclass axioms of the form $c_i \sqsubseteq c_j$, indicating that concept c_i is a subclass of c_j .
- **NF2.** Conjunctive subclass axioms $c_i \sqcap c_j \sqsubseteq c_k$, stating that the intersection of concepts c_i and c_j is a subclass of c_k .
- **NF3.** Existential inclusion axioms of the form $c_i \sqsubseteq \exists R.c_j$, meaning that instances of c_i are related via relation R to some instance of c_j .
- **NF4.** Existential restriction axioms $\exists R.c_i \sqsubseteq c_j$, expressing that any entity related to an instance of c_i via relation R must belong to concept c_j .

These normalized forms serve as the basis for encoding ontological constraints as geometric relations within the embedding space. The training loss is defined as

$$\mathcal{L}_{\text{axiom}} = \mathcal{L}_{\text{NF1}} + \mathcal{L}_{\text{NF2}} + \mathcal{L}_{\text{NF3}} + \mathcal{L}_{\text{NF4}}. \quad (14)$$

where L_1 to L_4 represent the training losses of the four axioms under ProtSAE. Appendix A.2 presents the detailed formulation of NF1 to NF4 and the derivation of the corresponding training loss.

4.3 Training Strategy

The overall training objective combines the reconstruction loss, the supervised prediction loss, and a semantic regularization term guided by domain knowledge:

$$\mathcal{L} = \underbrace{\|\hat{\mathbf{x}} - \mathbf{x}\|_2^2}_{\mathcal{L}_{\text{rec}}} + \lambda_{\text{annot}} \mathcal{L}_{\text{annot}} + \lambda_{\text{axiom}} \mathcal{L}_{\text{axiom}}, \quad (15)$$

where $\hat{\mathbf{x}}$ is the reconstructed input defined in Eq. (2). Here, \mathcal{L}_{rec} encourages faithful reconstruction of the input, $\mathcal{L}_{\text{annot}}$ is the cross-entropy loss in Eq. (7), $\mathcal{L}_{\text{axiom}}$ regularizes semantic alignment based on domain knowledge in Eq. (14). We use λ_{annot} and λ_{axiom} to weight these losses. Appendix B describes the detailed computation process.

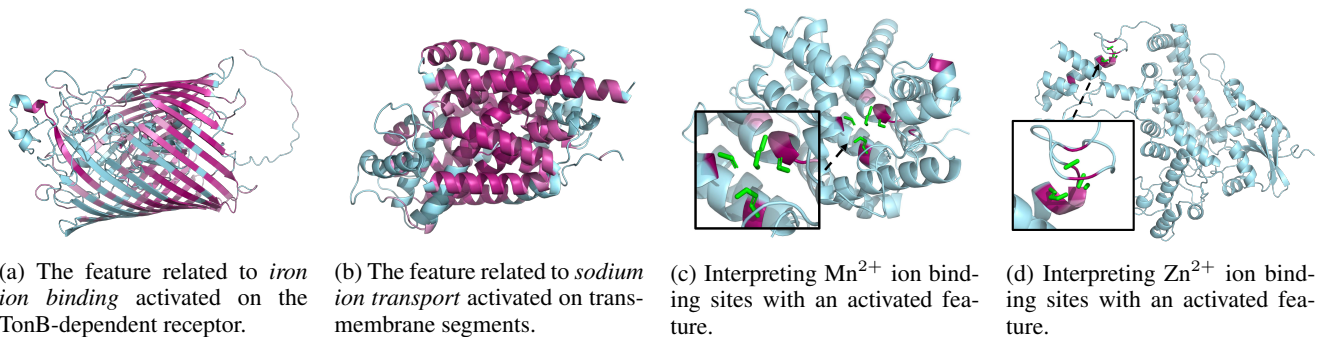


Figure 3: Interpretability visualization shows that ProtSAE reveals semantic alignment between learned features and protein structures, including functional regions and ion binding sites. We use red intensity to indicate feature activation strength, and green sticks to mark ground truth binding sites.

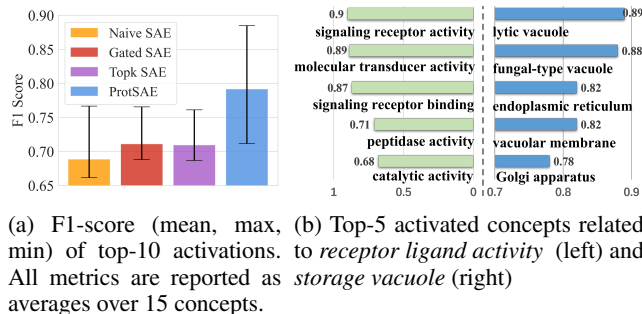


Figure 4: Comparison of relevance-based interpretation

5 Experiments and Results

5.1 Experiment Setup

Dataset and baselines. We construct the annotation data from protein function prediction datasets (Kulmanov et al. 2024). The protein function prediction datasets contains 77,647 proteins extracted from UniProtKB/Swiss-Prot. The annotated concepts can be categorized into three sub-ontologies: molecular function (MFO), biological process (BPO), and cellular component (CCO). Furthermore, we use the ion binding sites dataset (Yuan et al. 2022) as another annotation dataset, which covers four biologically relevant ion types: Zn^{2+} , Ca^{2+} , Mg^{2+} , and Mn^{2+} . We compare ProtSAE with widely adopted SAE baselines, including Naive SAE, Gated SAE (Rajamanoharan et al. 2024), and TopK SAE (Gao et al. 2024), in terms of both interpretability experiments and performance analyses. In the probing-based interpretation experiments, we further compare linear probing on the PLM hidden representations and the dictionary learning method SpLiCE (Bhalla et al. 2024). Detailed datasets and baseline settings are included in Appendix C.4.

Evaluation metrics. For relevance-based interpretation, we use F1-score to evaluate the relevance between neurons and biological concepts. To evaluate the probing-based interpretation, we employ standard metrics from the protein function prediction benchmark (Kulmanov et al. 2024) including AUPR, AUC, maximum protein-centric F-measure (F_{\max}),

and minimum semantic distance (S_{\min}). We use Loss Recovered (Rajamanoharan et al. 2024) to assess the reconstruction fidelity on varying sparsity. For intervention, we evaluate structural similarity using Template Modeling score (TM-score) (Zhang and Skolnick 2004) and Root Mean Square Distance (RMSD) (Betancourt and Skolnick 2001). Details are described in Appendix C.4.

Implementation. We train all SAE on the internal activations of ESM2-15B, with $5e^{-4}$ learning rate, 12,800 batch size, and 25,000 steps. For ProtSAE and TopK SAE, the number of active neurons K is varied in $\{50, 100, 500, 1000\}$. Gated SAE is tuned with L1 coefficients in $\{1.5e^{-4}, 2e^{-4}, 3e^{-4}, 4e^{-4}, 5e^{-4}\}$, and Naive SAE in $\{8e^{-5}, 6e^{-5}, 2e^{-4}, 3e^{-4}, 4e^{-4}\}$. All SAE activation width is set to 40,000 for BPO, 30,000 for MFO and CCO, and 10,000 for the ion binding-site dataset. λ_{annot} and λ_{axiom} are fixed at 1. Experiments are run on four NVIDIA A800 GPUs.

5.2 Interpretability Experiments

Interpretability visualization. Figure 3 visualizes the features learned by ProtSAE that are aligned with biological concepts and demonstrates their utility in interpreting and exploring the semantics of specific protein structural elements. Warmer colors (e.g., red) indicate stronger activation of the feature at that amino acid. The positions of the true binding sites are marked by green sticks. In Figure 3a, we examine activations associated with the concept *iron ion binding* on protein P06971. We observe strong activation in regions corresponding to the *TonB-dependent receptor* structure, which is known to be tightly associated with the recognition and transport of Fe^{3+} ions. In Figure 3b, the feature related to *sodium ion transport* is highly activated on the transmembrane segments of protein O67854, suggesting that certain α -helical transmembrane regions may play a crucial role in sodium ion transport. Furthermore, Figures 3c and 3d show that features related to *metal ion binding sites* can highlight binding sites in proteins.

Relevance-based interpretation evaluation. In this experiment, we evaluate whether ProtSAE can effectively identify features related to specific concepts. Following the previous work (Simon and Zou 2024), we extract relevant sequences

Method	$F_{\max} \uparrow$	$S_{\min} \downarrow$	AUPR \uparrow	AUC \uparrow
SpLiCE	.417	23.4	.360	.329
Naive SAE	.421	23.3	.340	.511
Gated SAE	.441	22.7	.368	.533
TopK SAE	.444	22.7	.379	.565
Linear Probe	.537	20.9	.522	.751
ProtSAE	.579	20.9	.487	.797

Table 1: Average performance across three datasets on probing-based interpretation

from UniProtKB,¹ and construct a validation set of 5,000 sequences for each of 15 GO terms along with an equal number of unrelated sequences used as negative examples. We compute the activation of each feature with respect to a given concept and report the top-10 features ranked by F1-score. A higher F1-score indicates a stronger correlation between the feature and the concept.

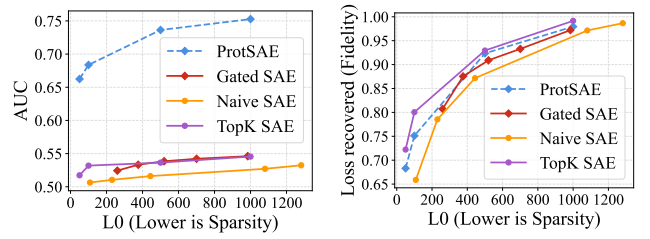
ProtSAE identifies features that are more semantically aligned with the target concepts. Figure 4a shows the average results of 15 concepts. Compared to the baselines, ProtSAE demonstrates significantly stronger relevance in both the mean and maximum activation levels. This suggests that the incorporation of semantic guidance during training enables ProtSAE to effectively disentangle semantic signals, and learns more accurate concept-related features. We further present additional highly activated features that show strong relevance to the validation set. As shown in Figure 4b, these features exhibit close functional or structural relationships with the target concept. For example, in the case of *storage vacuole*, the model activates structurally similar vacuole types such as *lytic vacuole* and *fungus-type vacuole*, as well as membrane-associated components like *vacuolar membrane* and *Golgi apparatus*.

Probing-based interpretation evaluation. To better evaluate the capability of probing-based interpretation, we conduct protein function prediction across datasets from three ontologies. The averaged results are summarized in Table 1. For a fair comparison, all SAE-based methods are evaluated under the same sparsity level. Notably, ProtSAE consistently outperforms all SAE baselines and the dictionary learning method SpLiCE across all evaluation metrics, and achieves comparable performance to linear probing on the hidden representations of PLMs. It suggests that ProtSAE, with semantic guidance, encourages the model to attend more effectively to the biological concepts during training. As a result, it mitigates the semantic loss that may occur during the SAE training, thereby achieving performance comparable to direct linear probing on PLMs.

5.3 Performance Analyses

Performance across different sparsity. Figure 5 illustrates the effect of sparsity on model performance. We assess the reconstruction fidelity using the metrics of *Loss Recovered*, which measures the proportion of the original PLM loss that can be recovered using SAE. The left subfigure shows the

¹<https://www.uniprot.org/>



(a) AUC under different sparsity (b) Loss Recovered under different sparsity

Figure 5: Performance comparison under different sparsity on the BPO dataset

AUC performance under varying levels of sparsity. ProtSAE consistently outperforms all baselines, indicating its ability to preserve semantics relevant to predefined concepts even under high sparsity, thereby achieving superior predictive performance. The right subfigure shows the trend of reconstruction fidelity as sparsity increases. Compared to other SAE variants, ProtSAE maintains comparable reconstruction quality, demonstrating its effectiveness in decomposing semantic concepts while faithfully preserving the original latent representations from PLMs. Detailed results on all datasets are provided in Appendix D.2.

Ablation study. We conduct an ablation study of ProtSAE by removing key components and retraining the model under various sparsity levels. Figure 6 depicts the ablation results. We discuss the effect of each component below:

- Without `detach`. In this variant, we no longer detach \mathbf{W}_{pred} when constructing \mathbf{W}_{def} , allowing the gradients from \mathcal{L}_{rec} to update \mathbf{W}_{pred} directly. So, the defined activations are now updated not only by the supervised data, but also by the reconstruction objective. While this slightly improves reconstruction fidelity, it leads to a dramatic decrease in AUC. This suggests that allowing \mathcal{L}_{rec} to influence \mathbf{W}_{pred} introduces entangled or ambiguous semantics into the defined activations, thereby degrading precision and interpretability.
- Removing $\mathcal{L}_{\text{axiom}}$. When the axiom learning component based on ELEmbeddings is removed, the training of \mathbf{W}_{pred} relies solely on the supervised data. This leads to a clear degradation in both AUC and reconstruction fidelity, highlighting the importance of modeling complex concept relationships through axioms to capture the intricate semantic structure of protein functions.
- Without \mathbf{z}_{bias} . In the right subfigure of Figure 6, we report the proportion of defined activations predicted as active that are indeed used during decoding. With \mathbf{z}_{bias} , nearly all predicted activations participate in reconstruction, indicating that \mathbf{z}_{def} holds strong potential for steering. Removing \mathbf{z}_{bias} reduces this proportion, weakening the alignment between prediction and actual activation.
- Without \mathbf{r}_{pred} . By setting the scaling parameter \mathbf{r}_{pred} in Eq. (8) to zero, \mathbf{W}_{def} and \mathbf{W}_{pred} become identical. This modification causes drops in both AUC and reconstruction

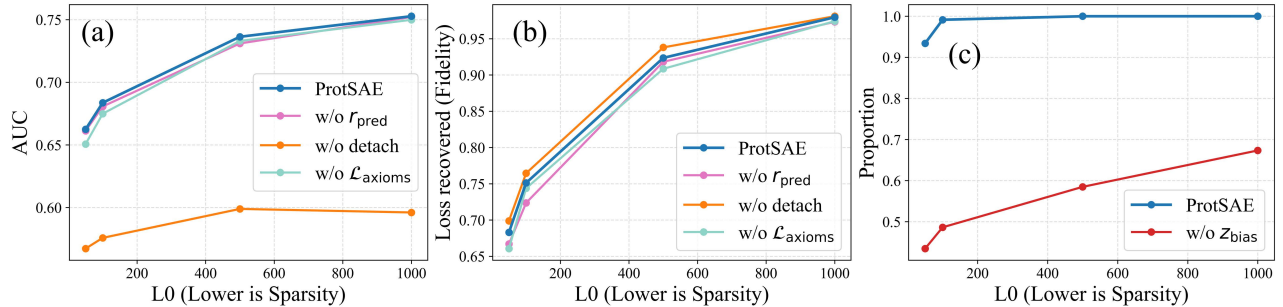


Figure 6: Ablation results on (a) AUC, (b) Loss Recovered, and (c) reconstruction proportion of predicted activations w.r.t. L_0

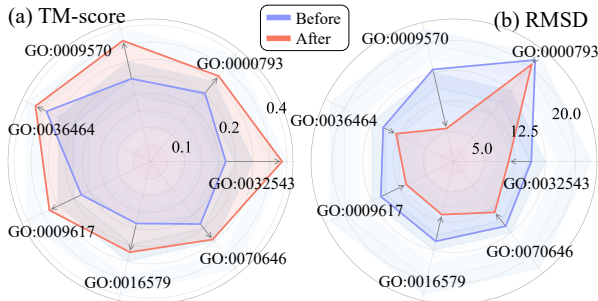


Figure 7: Effect of TM-score (left) and RMSD (right) before and after intervention

tion fidelity. The removal of r_{pred} hinders the model’s ability to learn feature magnitudes.

5.4 Steering Experiment

Concept intervention. We conduct a steering experiment across various biological concepts to evaluate whether ProtSAE can effectively steer PLM’s generation based on the learned concept-specific features. For each of seven selected concept-related sequences, we mask 50% of the tokens and compare the reconstructions generated before and after intervention. Following previous works (Zongying et al. 2024; Lv et al. 2024; Liu et al. 2025), we use TM-score and RMSD to measure structural similarity between the generated sequences and the natural proteins with the target concept, in order to assess whether the intervention can guide the model’s generation aligned with the desired concept. We use pLDDT to evaluate the structure stability.

As shown in Figure 7, TM-scores significantly increase while RMSD decreases after intervention, indicating improved structural alignment with the target concepts. Furthermore, Appendix E.2 includes detailed results and highlights significant improvements in the pLDDT scores of the generated proteins. These results suggest that ProtSAE successfully stores concept-aligned representations in its learned dictionary, and activating these features during generation enables the PLM to produce structurally stable proteins that better reflect the semantics of the desired concept.

Case study. Figure 8 visualizes proteins generated by ProtSAE after intervention (in blue). We identify their most

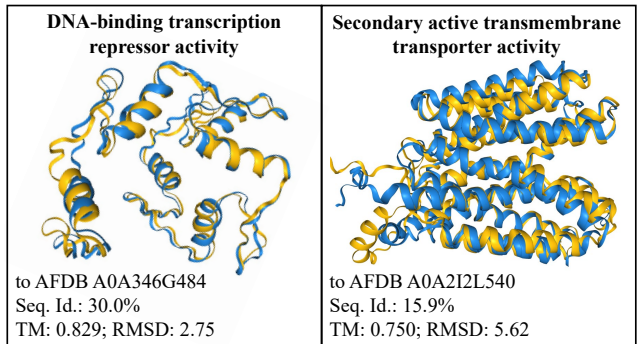


Figure 8: Intervention case study

similar natural counterparts (in yellow) with Foldseek (van Kempen et al. 2022). After intervention, ProtSAE can generate proteins with high structural similarity to natural counterparts with relevant concepts, remaining low sequence identity. This demonstrates that ProtSAE effectively captures concept-specific structural features and can successfully steer PLM’s generation. For example, we intervene on the concept of “DNA-binding transcription repressor activity” and generate a protein structurally similar to the natural protein “A0A346G484” (TM-score: 0.829, RMSD:2.75), while maintaining sequence novelty (Seq. ID: 30.0%). “A0A346G484” contains a putative zinc-finger domain and is annotated with the desired concept. We also generate proteins with high structural similarity to natural proteins exhibiting transmembrane transporter activity.

6 Conclusion

We propose ProtSAE, a semantically-guided SAE to tackle semantic entanglement in SAE training and improve interpretability of PLMs. We introduce domain knowledge into ProtSAE to constrain the relationship among concepts, and apply forced activations and feature rescaling to ensure that the learned features effectively contribute to the reconstruction while maintaining high reconstruction fidelity. Interpretability experiments show that ProtSAE consistently captures features more aligned with protein structures and functions. Performance analyses and steering experiments show the superiority of ProtSAE against existing SAE baselines.

References

- Adams, E.; Bai, L.; Lee, M.; Yu, Y.; and AlQuraishi, M. 2025. From Mechanistic Interpretability to Mechanistic Biology: Training, Evaluating, and Interpreting Sparse Autoencoders on Protein Language Models. *bioRxiv*, 2025–02.
- Betancourt, M. R.; and Skolnick, J. 2001. Universal similarity measure for comparing protein structures. *Biopolymers: Original Research on Biomolecules*, 59(5): 305–309.
- Bhalla, U.; Oesterling, A.; Srinivas, S.; Calmon, F.; and Lakkaraju, H. 2024. Interpreting clip with sparse linear concept embeddings (splice). *NeurIPS*, 37: 84298–84328.
- Clark, W. T.; and Radivojac, P. 2013. Information-theoretic evaluation of predicted ontological annotations. *Bioinformatics*, 29(13): i53–i61.
- Dong, T.; Kan, C.; Devkota, K.; and Singh, R. 2024. Allo-Allo: Data-efficient prediction of allosteric sites. *bioRxiv*, 2024–09.
- Dunefsky, J.; Chlenski, P.; and Nanda, N. 2024. Transcoders find interpretable LLM feature circuits. In *NeurIPS*.
- Elhage, N.; Hume, T.; Olsson, C.; Schiefer, N.; Henighan, T.; Kravec, S.; Hatfield-Dodds, Z.; Lasenby, R.; Drain, D.; Chen, C.; et al. 2022. Toy models of superposition. *arXiv*.
- Ferruz, N.; and Höcker, B. 2022. Controllable protein design with language models. *Nature Machine Intelligence*, 4(6): 521–532.
- Gallifant, J.; Chen, S.; Sasse, K.; Aerts, H.; Hartvigsen, T.; and Bitterman, D. S. 2025. Sparse autoencoder features for classifications and transferability. *arXiv*.
- Gao, L.; la Tour, T. D.; Tillman, H.; Goh, G.; Troll, R.; Radford, A.; Sutskever, I.; Leike, J.; and Wu, J. 2024. Scaling and evaluating sparse autoencoders. *arXiv*.
- Garcia, E. N. V.; and Ansuini, A. 2025. Interpreting and Steering Protein Language Models through Sparse Autoencoders. *arXiv*.
- Gurnee, W.; Nanda, N.; Pauly, M.; Harvey, K.; Troitskii, D.; and Bertsimas, D. 2023. Finding Neurons in a Haystack: Case Studies with Sparse Probing. *Transactions on Machine Learning Research*.
- Ismail, A. A.; Oikarinen, T.; Wang, A.; Adebayo, J.; Stanton, S. D.; Bravo, H. C.; Cho, K.; and Frey, N. C. 2025. Concept Bottleneck Language Models For Protein Design. In *ICLR*.
- Joshi, S.; Dittadi, A.; Lachapelle, S.; and Sridhar, D. 2025. Identifiable Steering via Sparse Autoencoding of Multi-Concept Shifts. *arXiv*.
- Kannan, G. R.; Hie, B. L.; and Kim, P. S. 2024. Single-Sequence, Structure Free Allosteric Residue Prediction with Protein Language Models. *bioRxiv*, 2024–10.
- Kulmanov, M.; Guzmán-Vega, F. J.; Duek Roggli, P.; Lane, L.; Arold, S. T.; and Hoehndorf, R. 2024. Protein function prediction as approximate semantic entailment. *Nature Machine Intelligence*, 6(2): 220–228.
- Kulmanov, M.; and Hoehndorf, R. 2022. DeepGOZero: Improving protein function prediction from sequence and zero-shot learning based on ontology axioms. *Bioinformatics*, 38(Supplement_1): i238–i245.
- Kulmanov, M.; Liu-Wei, W.; Yan, Y.; and Hoehndorf, R. 2019. EL Embeddings: Geometric construction of models for the description logic EL++. In *IJCAI*, 6103–6109.
- Lin, B.; Luo, X.; Liu, Y.; and Jin, X. 2024. A comprehensive review and comparison of existing computational methods for protein function prediction. *Briefings in Bioinformatics*, 25(4): 289.
- Lin, Z.; Akin, H.; Rao, R.; Hie, B.; Zhu, Z.; Lu, W.; Smetanin, N.; Verkuil, R.; Kabeli, O.; Shmueli, Y.; et al. 2023a. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637): 1123–1130.
- Lin, Z.; Akin, H.; Rao, R.; Hie, B.; Zhu, Z.; Lu, W.; Smetanin, N.; Verkuil, R.; Kabeli, O.; Shmueli, Y.; et al. 2023b. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637): 1123–1130.
- Liu, X.; Liu, Y.; Chen, S.; and Hu, W. 2025. Controllable Protein Sequence Generation with LLM Preference Optimization. In *AAAI*.
- Lv, L.; Lin, Z.; Li, H.; Liu, Y.; Cui, J.; Chen, C. Y.-C.; Yuan, L.; and Tian, Y. 2024. ProLLaMA: A protein large language model for multi-task protein language processing. *arXiv preprint arXiv:2402.16445*.
- Makelov, A. 2024. Sparse autoencoders match supervised features for model steering on the ioi task. In *ICML 2024 Workshop on Mechanistic Interpretability*.
- Makelov, A.; Lange, G.; Geiger, A.; and Nanda, N. 2024. Is This the Subspace You Are Looking for? An Interpretability Illusion for Subspace Activation Patching. In *ICLR*.
- Miller, J.; Chughtai, B.; and Saunders, W. 2024. Transformer Circuit Faithfulness Metrics are not Robust. In *COLM*.
- Nijkamp, E.; Ruffolo, J. A.; Weinstein, E. N.; Naik, N.; and Madani, A. 2023. ProGen2: Exploring the boundaries of protein language models. *Cell Systems*, 14(11): 968–978.
- Nori, D.; Singireddy, S.; and Have, M. T. 2023a. Identification of Knowledge Neurons in Protein Language Models. *arXiv preprint arXiv:2312.10770*.
- Nori, D.; Singireddy, S.; and Have, M. T. 2023b. Identification of Knowledge Neurons in Protein Language Models. *arXiv preprint arXiv:2312.10770*.
- Olah, C.; Cammarata, N.; Schubert, L.; Goh, G.; Petrov, M.; and Carter, S. 2020. Zoom in: An introduction to circuits. *Distill*, 5(3): e00024–001.
- Pach, M.; Karthik, S.; Bouniot, Q.; Belongie, S.; and Akata, Z. 2025. Sparse Autoencoders Learn Monosemantic Features in Vision-Language Models. *arXiv preprint arXiv:2504.02821*.
- Paulo, G.; Mallen, A.; Juang, C.; and Belrose, N. 2024. Automatically interpreting millions of features in large language models. *arXiv preprint arXiv:2410.13928*.
- Radivojac, P.; Clark, W. T.; Oron, T. R.; Schnoes, A. M.; Wittkop, T.; Sokolov, A.; Graim, K.; Funk, C.; Verspoor, K.; Ben-Hur, A.; et al. 2013. A large-scale evaluation of computational protein function prediction. *Nature Methods*, 10(3): 221–227.

Rajamanoharan, S.; Conmy, A.; Smith, L.; Lieberum, T.; Varma, V.; Kramar, J.; Shah, R.; and Nanda, N. 2024. Improving sparse decomposition of language model activations with gated sparse autoencoders. In *NeurIPS*.

Rao, R.; Meier, J.; Sercu, T.; Ovchinnikov, S.; and Rives, A. 2021. Transformer protein language models are unsupervised structure learners. In *ICLR*.

Shi, C.; Beltran Velez, N.; Nazaret, A.; Zheng, C.; Garriga-Alonso, A.; Jesson, A.; Makar, M.; and Blei, D. 2024. Hypothesis testing the circuit hypothesis in LLMs. In *NeurIPS*.

Simon, E.; and Zou, J. 2024. InterPLM: Discovering Interpretable Features in Protein Language Models via Sparse Autoencoders. *bioRxiv*, 2024–11.

Tamkin, A.; Tafeeque, M.; and Goodman, N. 2024. Codebook Features: Sparse and Discrete Interpretability for Neural Networks. In *ICML*, 47535–47563.

Todd, E.; Li, M.; Sharma, A. S.; Mueller, A.; Wallace, B. C.; and Bau, D. 2024. Function Vectors in Large Language Models. In *ICLR*.

van Kempen, M.; Kim, S. S.; Tumescheit, C.; Mirdita, M.; Gilchrist, C. L.; Söding, J.; and Steinegger, M. 2022. Foldseek: Fast and accurate protein structure search. *bioRxiv*, 2022–02.

Vig, J.; Madani, A.; Varshney, L. R.; Xiong, C.; Rajani, N.; et al. 2021. BERTology Meets Biology: Interpreting Attention in Protein Language Models. In *ICLR*.

Wang, K. R.; Variengien, A.; Conmy, A.; Shlegeris, B.; and Steinhardt, J. 2023. Interpretability in the Wild: A Circuit for Indirect Object Identification in GPT-2 Small. In *ICLR*.

Wright, B.; and Sharkey, L. 2024. Addressing feature suppression in saes. In *AI Alignment Forum*, 16.

Wu, X.; Yu, W.; Zhai, X.; and Liu, N. 2025a. Self-regularization with latent space explanations for controllable LLM-based classification. *arXiv*.

Wu, X.; Yuan, J.; Yao, W.; Zhai, X.; and Liu, N. 2025b. Interpreting and steering LLMs with mutual information-based explanations on sparse autoencoders. *arXiv*.

Yuan, Q.; Chen, S.; Wang, Y.; Zhao, H.; and Yang, Y. 2022. Alignment-free metal ion-binding site prediction from protein sequence through pretrained language model and multi-task learning. *Briefings in Bioinformatics*, 23(6): bbac444.

Yun, Z.; Chen, Y.; Olshausen, B. A.; and LeCun, Y. 2021. Transformer visualization via dictionary learning: Contextualized embedding as a linear superposition of transformer factors. *NAACL-HLT*, 1.

Zhang, Y.; and Skolnick, J. 2004. Scoring function for automated assessment of protein structure template quality. *Proteins: Structure, Function, and Bioinformatics*, 57(4): 702–710.

Zhang, Z.; Wayment-Steele, H. K.; Brixi, G.; Wang, H.; Kern, D.; and Ovchinnikov, S. 2024. Protein language models learn evolutionary statistics of interacting sequence motifs. *Proceedings of the National Academy of Sciences*, 121(45): e2406285121.

Zongying, L.; Hao, L.; Liuzhenghao, L.; Bin, L.; Junwu, Z.; Yu-Chian, C. C.; Li, Y.; and Yonghong, T. 2024. TaxDiff: Taxonomic-Guided Diffusion Model for Protein Sequence Generation. *arXiv preprint arXiv:2402.17156*.

A Derivation of ELEmbeddings in SAE

A.1 Structural Equivalence Between ELEmbeddings Prediction Function and SAE Encoder

We formally establish the structural equivalence between the ELEmbeddings prediction function and the forward pass of the encoder in SAE.

In the ELEmbeddings framework, given a protein sequence p and a concept class c_i , the prediction score is defined as

$$y'_i = \sigma \left(f_\eta(p)^\top \cdot (f_\eta(hF) + f_\eta(c_i)) + r_\eta(c_i) \right), \quad (16)$$

where $f_\eta(\cdot)$ denotes a projection into the ontology embedding space of dimension d , hF represents the hierarchical context of the ontology, and $r_\eta(c_i) > 0$ is a learnable bias term corresponding to a soft radius threshold. In our setting, protein functions are inferred by performing token-level semantic predictions, where each token representation is evaluated w.r.t. specific biological functions. The protein p is first tokenized as $p = \{t_1, t_2, \dots, t_L\}$, and each token t_j is embedded using a pretrained language model. Let $\mathbf{x}_j \in \mathbb{R}^d$ denote the hidden representation of token t_j , which we treat as the input to the encoder. That is, we interpret

$$f_\eta(t_j) = \mathbf{x}_j - \mathbf{b}_{\text{dec}}, \quad (17)$$

and define

$$\mathbf{w}_i := f_\eta(hF) + f_\eta(c_i) \in \mathbb{R}^d, b_i := r_\eta(c_i) \in \mathbb{R}_{>0}. \quad (18)$$

Then, Eq. (16) becomes, for each token j ,

$$y'_{i,j} = \sigma \left(\mathbf{w}_i^\top (\mathbf{x}_j - \mathbf{b}_{\text{dec}}) + b_i \right), \quad (19)$$

which mirrors the token-level prediction of a sparse encoder.

In ProtSAE, each concept c_i is predicted per token as

$$\pi_{\text{pred}}^{(i,j)} = \sigma \left(\mathbf{W}_{\text{pred}}^{(i)} (\mathbf{x}_j - \mathbf{b}_{\text{dec}}) + b_{\text{pred}}^{(i)} \right), \quad (20)$$

where $\mathbf{W}_{\text{pred}}^{(i)} \in \mathbb{R}^d$ and $b_{\text{pred}}^{(i)} \in \mathbb{R}$ are learned parameters.

By assigning

$$\mathbf{W}_{\text{pred}}^{(i)} := \mathbf{w}_i^\top = f_\eta(hF) + f_\eta(c_i), \quad (21)$$

$$b_{\text{pred}}^{(i)} := b_i = r_\eta(c_i), \quad \text{with } b_i > 0 \quad (22)$$

we obtain

$$\pi_{\text{pred}}^{(i,j)} = y'_{i,j}, \quad (23)$$

which demonstrates the structural equivalence of Eq. (20) and Eq. (16) at the token level.

To compute the sequence-level concept prediction $\pi_{\text{pred}}^{(i)}$ for the full protein sequence, token-level predictions are aggregated (e.g., via max-pooling or average-pooling) across

all positions, which is similar to the interpretability procedure in traditional SAE:

$$\pi_{\text{pred}}^{(i)} = \text{Pool}_j \left(\pi_{\text{pred}}^{(i,j)} \right). \quad (24)$$

From Eq. (22), we establish a direct correspondence between the encoder weights of the SAE and the terms embedding defined in ELEmbeddings. Under the constraint $b_i > 0$, which ensures a valid geometric interpretation of b_i as the radius of a high-dimensional sphere, the encoder weight vector $\mathbf{W}_{\text{pred}}^{(i)}$ learned by the SAE can be interpreted as the sum of two components: the embedding of the ontology concept c_i , and the offset vector representing the has-Function relation $f_\eta(hF)$. Therefore, the direction of each encoder unit in the SAE corresponds to the logical composition $f_\eta(hF) + f_\eta(c_i)$ in the ELEmbeddings framework.

This structural equivalence implies that the ontology-aware semantic geometry encoded by ELEmbeddings is preserved in the SAE encoder. As a result, the inter-term relations of the concepts, such as subsumption or regulation, can be explicitly modeled by analyzing the encoder matrix \mathbf{W}_{pred} . This provides a principled way to ground LLM activations in interpretable, structured biological semantics.

A.2 Training SAE with ELEmbeddings Axioms

ELEmbeddings utilize four normalized axiom forms (NF1 to NF4) to encode ontological constraints. For each normalized form, a specific geometric loss function is defined. Below, we present the detailed loss formulations and explanations.

From Eq. (22), we observe that the ontology-based representation $f_d(c_i)$ for a given concept c_i can be directly computed using the encoder weight $\mathbf{W}_{\text{pred}}^{(i)}$ obtained from the SAE. This correspondence enables us to directly compute the ELEmbeddings losses with the encoder weights, without requiring separate embedding training. For all involved relations, we initialize them in the same way as in \mathbf{W}_{pred} .

NF1 Loss. NF1 corresponds to simple subclass axioms of the form $c_i \sqsubseteq c_j$, e.g., “binding” (GO:0005488) SubClassOf “molecular function” (GO:0003674). The corresponding loss penalizes the distance between the centers of the two n -balls and ensures that the n -ball for c_i is fully contained in the n -ball for c_j :

$$\begin{aligned} \mathcal{L}_{\text{NF1}} &= \frac{1}{|\text{NF1}|} \sum_{c_i, c_j \in \text{NF1}} \max(0, \|f_\eta(c_i) - f_\eta(c_j)\| \\ &\quad + r_\eta(c_i) - r_\eta(c_j) - \gamma) \\ &= \frac{1}{|\text{NF1}|} \sum_{i,j \in \text{NF1}} \max(0, \|\mathbf{w}_i - \mathbf{w}_j\| + b_i - b_j - \gamma). \end{aligned} \quad (25)$$

NF2 Loss. NF2 handles axioms of the form $c_i \sqcap c_j \sqsubseteq c_k$, which express that the intersection of two concepts is a subclass of a third, e.g., “cutinase activity” (GO:0050525) and “biological regulation” (GO:0065007) SubClassOf “positive regulation of protein kinase B signaling” (GO:0051897). The corresponding loss minimizes

the discrepancy between the intersection of n -balls for c_i and c_j , and the n -ball for class E :

$$\begin{aligned} \mathcal{L}_{\text{NF2}} &= \frac{1}{|\text{NF2}|} \sum_{c_i, c_j, c_k \in \text{NF2}} \left(\max(0, \|f_\eta(c_i) - f_\eta(c_j)\| \right. \\ &\quad \left. - r_\eta(c_i) - r_\eta(c_j) - \gamma) \right. \\ &\quad \left. + \max(0, \|f_\eta(c_i) - f_\eta(c_k)\| - r_\eta(c_i) - \gamma) \right. \\ &\quad \left. + \max(0, \|f_\eta(c_j) - f_\eta(c_k)\| - r_\eta(c_j) - \gamma) \right. \\ &\quad \left. + \max(0, \min(r_\eta(c_i), r_\eta(c_j)) - r_\eta(c_k) - \gamma) \right) \\ &= \frac{1}{|\text{NF2}|} \sum_{i,j,k \in \text{NF2}} \left(\max(0, \|\mathbf{w}_i - \mathbf{w}_j\| - b_i - b_j - \gamma) \right. \\ &\quad \left. + \max(0, \|\mathbf{w}_i - \mathbf{w}_k\| - b_i - \gamma) \right. \\ &\quad \left. + \max(0, \|\mathbf{w}_j - \mathbf{w}_k\| - b_j - \gamma) \right. \\ &\quad \left. + \max(0, \min(b_i, b_j) - b_k - \gamma) \right). \end{aligned} \quad (26)$$

NF3 Loss. NF3 corresponds to axioms of the form $c_i \sqsubseteq \exists R.c_j$, indicating that concept c_i is included in the set of entities that are related by R to some instance of c_j , e.g., “positive regulation of arginine biosynthetic process” (GO:1900080) SubClassOf “positively regulates” (RO:0002213) some “arginine biosynthetic process” (GO:0006526). This is modeled by translating the n -ball of class c_j by the relation vector $f_\eta(R)$ and minimizing its non-overlap with c_i :

$$\begin{aligned} \mathcal{L}_{\text{NF3}} &= \frac{1}{|\text{NF3}|} \sum_{R, c_i, c_j \in \text{NF3}} \max(0, \|f_\eta(c_i) - f_\eta(R) - f_\eta(c_j)\| \\ &\quad - r_\eta(c_i) - r_\eta(c_j) - \gamma) \\ &= \frac{1}{|\text{NF3}|} \sum_{R, c_i, c_j \in \text{NF3}} \max(0, \|\mathbf{w}_i - \mathbf{w}_j - f_\eta(R)\| \\ &\quad - b_i - b_j - \gamma). \end{aligned} \quad (27)$$

NF4 Loss. NF4 axioms are of the form $\exists R.c_i \sqsubseteq c_j$, implying that entities related by R to some instance of c_i are contained in class c_j , e.g., “part of” (BFO:0000050) some “conjugation” (GO:0000746) SubClassOf “mammary stem cell proliferation” (GO:0002174). The loss translates the n -ball of c_i by relation vector $f_\eta(R)$, and ensures containment within c_j ’s n -ball:

$$\begin{aligned} \mathcal{L}_{\text{NF4}} &= \frac{1}{|\text{NF4}|} \sum_{c_i, R, c_j \in \text{NF4}} \max(0, \|f_\eta(c_i) + f_\eta(R) - f_\eta(c_j)\| \\ &\quad + r_\eta(c_i) - r_\eta(c_j) - \gamma) \\ &= \frac{1}{|\text{NF4}|} \sum_{c_i, R, c_j \in \text{NF4}} \max(0, \|\mathbf{w}_i - \mathbf{w}_j + f_\eta(R)\| \\ &\quad + b_i - b_j - \gamma). \end{aligned} \quad (28)$$

Algorithm 1: Forward Pass with SAE

Input: Input $\mathbf{x} \in \mathbb{R}^d$
Output: Reconstruction $\hat{\mathbf{x}}$, semantic prediction π_{pred}

- 1 $\mathbf{W}_{\text{def}} \leftarrow \mathbf{W}_{\text{pred}}^{\text{detach}} \cdot \exp(\mathbf{r}_{\text{pred}})$
/* Eq. (8), compute weight matrix */
- 2 $\mathbf{z}_{\text{unk}} \leftarrow \text{TopK}(\mathbf{W}_{\text{unk}}(\mathbf{x} - \mathbf{b}_{\text{dec}}) + \mathbf{b}_{\text{unk}})$
/* Eq. (10), compute activations */
- 3 $\hat{\mathbf{z}}_{\text{def}} \leftarrow \mathbf{W}_{\text{def}}(\mathbf{x} - \mathbf{b}_{\text{dec}}) + \mathbf{b}_{\text{def}}$
- 4 $\mathbf{z}_{\text{bias}} \leftarrow \mathbb{1}_{\pi_{\text{pred}} > 0} \cdot \text{ReLU}(\text{mean}(\mathbf{z}_{\text{unk}}) - \hat{\mathbf{z}}_{\text{def}})$
/* Eq. (12), force activation */
- 5 $\mathbf{z}_{\text{def}} \leftarrow \hat{\mathbf{z}}_{\text{def}} + \mathbf{z}_{\text{bias}}$
/* Eq. (12), defined activations */
- 6 $\mathbf{z} \leftarrow \text{Concat}(\mathbf{z}_{\text{def}}, \mathbf{z}_{\text{unk}})$
- 7 $\hat{\mathbf{x}} \leftarrow \mathbf{W}_{\text{dec}}\mathbf{z} + \mathbf{b}_{\text{dec}}$
/* Eq. (2), reconstruction */
- 8 $\pi_{\text{pred}} \leftarrow \sigma(\mathbf{W}_{\text{pred}}(\mathbf{x} - \mathbf{b}_{\text{dec}}) + \mathbf{b}_{\text{pred}})$
/* Eq. (6), prediction */
- 9 **return** $\hat{\mathbf{x}}, \pi_{\text{pred}}$

Algorithm 2: Training Procedure

Input: Dataset $\mathcal{D} = \{(\mathbf{x}^{(j)}, y^{(j)})\}$ with semantic labels $y^{(j)} \in \{0, 1\}^m$

- 1 **foreach** *mini-batch* $\{\mathbf{x}, y\} \subset \mathcal{D}$ **do**
- 2 Compute $\mathbf{z}_{\text{def}}, \mathbf{z}_{\text{unk}}, \pi_{\text{pred}}, \hat{\mathbf{x}}$ using Algorithm 1
- 3 $\mathcal{L}_{\text{rec}} \leftarrow \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$
/* Reconstruction loss */
- 4 $\mathcal{L}_{\text{annot}} \leftarrow \text{CrossEntropy}(\pi_{\text{pred}}, y)$
/* Annotation guidance loss */
- 5 $\mathcal{L}_{\text{axioms}} \leftarrow \mathcal{L}_{\text{NF1}} + \mathcal{L}_{\text{NF2}} + \mathcal{L}_{\text{NF3}} + \mathcal{L}_{\text{NF4}}$
/* Get normalized axiom loss */
- 6 $\mathcal{L} \leftarrow \mathcal{L}_{\text{rec}} + \lambda_{\text{annot}}\mathcal{L}_{\text{annot}} + \lambda_{\text{axiom}}\mathcal{L}_{\text{axiom}}$
/* Total loss */
- 7 Backpropagate \mathcal{L} and update all trainable parameters

B Pseudocode

C Interpretability Experiments Appendix

C.1 Settings

To better demonstrate the superior interpretability of ProtSAE, we design comprehensive interpretability experiments, including both relevance-based and probing-based evaluations. Furthermore, we present several cases on ProtSAE interpretability including protein functional structure and metal ion binding site prediction.

For probing-based interpretation, we perform experiments on protein function prediction tasks. This task aims to uncover the biological roles and interactions of proteins, which is critical for applications such as drug target identification, understanding disease mechanisms, and advancing biotechnology. Protein annotations are sourced from the Gene Ontology (GO), which is structured into three sub-ontologies: Molecular Function (MFO), Biological Process (BPO), and Cellular Component (CCO).

For relevance-based interpretation, we analyze 15 biologically meaningful concepts drawn from three distinct datasets. For each concept, we retrieve relevant proteins from UniProtKB, ensuring that all training samples are excluded to construct an independent evaluation set consisting of 5,000 proteins. In this analysis, the top-10 most activated features identified by the model are selected as candidate features for interpretation.

C.2 Datasets

We adopt the protein function prediction benchmark from the previous work (Kulmanov et al. 2024).² The dataset is extracted from UniProtKB/Swiss-Prot and is filtered to retain those with experimental annotations supported by evidence codes such as EXP, IDA, IPI, IMP, IGI, IEP, TAS, IC, HTP, HDA, HMP, HGI, and HEP. This filtering resulted in a high-quality dataset containing 77,647 manually curated proteins. The Gene Ontology release used corresponds to November 16, 2021. Models are trained and evaluated separately for each GO sub-ontology. Detailed information is provided in Table 2. It shows the number of GO terms, total number of proteins, number of groups of similar proteins, number of proteins in training, validation and testing sets.

To further assess the interpretability of ProtSAE, we conduct experiments on the metal ion-binding datasets (Yuan et al. 2022). These datasets cover four biologically relevant ion types: Zn^{2+} , Ca^{2+} , Mg^{2+} , and Mn^{2+} . Detailed statistics for each subset are provided in Table 3.

C.3 Implementations

We train all SAE models on protein function prediction datasets, using a learning rate of $5e-4$ and a batch size of 12,800 for 25,000 steps. For ProtSAE and TopK SAE, we vary the number of active neurons with $K \in \{50, 100, 500, 1000\}$. For Gated SAE, we tune the L1 regularization coefficient in $\{1.5e-4, 2e-4, 3e-4, 4e-4, 5e-4\}$. For Naive SAE, we explore L1 coefficients in $\{8e-5, 6e-5, 2e-4, 3e-4, 4e-4\}$. We set the activation width to 40,000 for the BPO dataset and 30,000 for both MFO and CCO. The dimensionality of the activation vectors matches that of the ESM2-15B model’s hidden states (5,120). We fix both λ_{sup} and λ_{axiom} to 1. All results are reported using the representations from layer 35 of ESM2-15B. We train ProtSAE on the metal-ion binding dataset using a learning rate of 5×10^{-4} and a batch size of 3200 for 40,000 optimization steps. The activation dimension of ProtSAE is set to 10,000. For evaluation, we compute top- K activations with $K \in \{50, 100, 500, 1000\}$ to analyze feature sparsity and relevance. All experiments are conducted on four NVIDIA A800 GPUs.

C.4 Baselines and Metrics

We compare ProtSAE with several representative SAE baselines, a dictionary learning method SpLiCE (Bhalla et al. 2024) and linear probe on PLMs hidden representations:

²The dataset is publicly available under the BSD 3-Clause License.

Ontology	GO Terms	Proteins	Groups	Training	Validation	Testing
MFO	6,851	43,279	6,963	52,072	2,964	4,221
BPO	21,356	58,729	9,463	52,584	2,870	3,275
CCO	2,829	59,257	10,019	48,318	4,970	5,969

Table 2: Summary of the UniProtKB/Swiss-Prot dataset

Ligand type	Dataset	Binding residue	Non-binding residue
Zn^{2+}	ZN_Train_1647	7,731	467,184
	ZN_Test_211	1,039	54,981
Ca^{2+}	CA_Train_1554	8,442	495,700
	CA_Test_183	1,034	65,820
Mg^{2+}	MG_Train_1730	6,321	569,572
	MG_Test_235	893	87,913
Mn^{2+}	MN_Train_547	2,556	179,143
	MN_Test_57	225	20,194

Table 3: Statistics of the metal ion binding-sites prediction dataset

- **Naive SAE:** A standard SAE trained only with reconstruction loss and a simple L1 penalty for sparsity.
- **Gated SAE:** An enhanced SAE variant that applies learned gating mechanisms and L1 regularization to encourage selective activation.
- **TopK SAE:** A Top-K SAE that enforces hard sparsity by retaining only the top-K highest activations for each input.
- **SpLiCE:** A dictionary learning method that transforms representations into sparse linear combinations of human interpretable concepts.
- **Linear Probe:** We direct use linear probe on LLMs’ internal representation, to show the information loss of SAE methods.

Metrics. For relevance-based interpretation, following (Simon and Zou 2024; Garcia and Ansuini 2025), we use the F1 score to evaluate the relevance of selected features. The detailed calculation is as follows:

$$\text{precision} = \frac{\text{TruePositives}}{\text{TruePositives} + \text{FalsePositives}}, \quad (29)$$

$$\text{recall} = \frac{\text{TruePositive}}{\text{TruePositives} + \text{FalseNegative}}, \quad (30)$$

$$F1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}. \quad (31)$$

For each feature, we determine whether it is positive based on whether its normalized activation value exceeded a specific threshold. For each concept, we searched thresholds from the set $\{0, 0.05, 0.1, 0.15, 0.2, 0.25, 0.5, 0.8\}$, and selected the threshold that yielded the highest F1 score.

For probing-based prediction, we adopt AUPR and a class-centric AUC from (Kulmanov and Hoehndorf 2022; Kulmanov et al. 2024). We use the protein-centric evaluation metrics F_{\max} and S_{\min} introduced by the CAFA challenge (Clark and Radivojac 2013; Radivojac et al. 2013). The detailed formulations of metrics are provided below.

F_{\max} is a maximum protein-centric F-measure computed over all prediction thresholds. First, it computes average precision and recall using the following equations:

$$pr_u(\tau) = \frac{\sum_c I(c \in \hat{T}_u(\tau) \wedge c \in T_u)}{\sum_c I(c \in \hat{T}_u(\tau))}, \quad (32)$$

$$rc_u(\tau) = \frac{\sum_c I(c \in \hat{T}_u(\tau) \wedge c \in T_u)}{\sum_c I(c \in T_u)}, \quad (33)$$

$$\text{AvgPr}(\tau) = \frac{1}{|\text{set}(\tau)|} \sum_{u \in \text{set}(\tau)} pr_u(\tau), \quad (34)$$

$$\text{AvgRc}(\tau) = \frac{1}{N} \sum_{u=1}^N rc_u(\tau), \quad (35)$$

where c is a GO class, T_u is the set of true annotations, $\hat{T}_u(\tau)$ is the set of predicted annotations for protein u at threshold τ , $\text{set}(\tau)$ is the set of proteins for which at least one GO class is predicted at threshold τ , N is the total number of proteins, and I is the indicator function returning 1 if the condition holds and 0 otherwise. Then, we compute F_{\max} over thresholds $\tau \in [0, 1]$ with a step size of 0.01. A class c is considered predicted for protein u if its score is greater than or equal to τ :

$$F_{\max} = \max_{\tau} \left\{ \frac{2 \cdot \text{AvgPr}(\tau) \cdot \text{AvgRc}(\tau)}{\text{AvgPr}(\tau) + \text{AvgRc}(\tau)} \right\}. \quad (36)$$

S_{\min} computes the semantic distance between real and predicted annotations based on information content of the classes. The information content $IC(c)$ is computed based on the annotation probability of class c :

$$IC(c) = -\log(\text{Pr}(c | \mathcal{P}(c))), \quad (37)$$

where $\mathcal{P}(c)$ denotes the parent classes of c . The S_{\min} is computed using the following equation:

$$S_{\min} = \min_{\tau} \sqrt{ru(\tau)^2 + mi(\tau)^2}, \quad (38)$$

where $ru(\tau)$ is the average remaining uncertainty and $mi(\tau)$ is the average misinformation:

$$ru(\tau) = \frac{1}{N} \sum_{u=1}^N \sum_{c \in T_u \setminus \hat{T}_u(\tau)} IC(c), \quad (39)$$

$$mi(\tau) = \frac{1}{N} \sum_{u=1}^N \sum_{c \in \hat{T}_u(\tau) \setminus T_u} IC(c). \quad (40)$$

C.5 Metal ion binding-sites prediction cases.

In this section, we present a case study demonstrating the interpretability of ion binding-sites concepts. We use the normalized activation score of the feature corresponding to a specific ion binding-sites concept as the prediction signal from ProtSAE for identifying whether an amino acid residue is part of an ion binding site. In Figure 9, regions with higher activation scores are shown in red, while blue spheres indicate the locations of ground-truth binding sites. As observed, the learned feature is highly activated at the true binding sites while remaining largely inactive in non-target regions, suggesting a strong correlation between the feature and the intended concept.

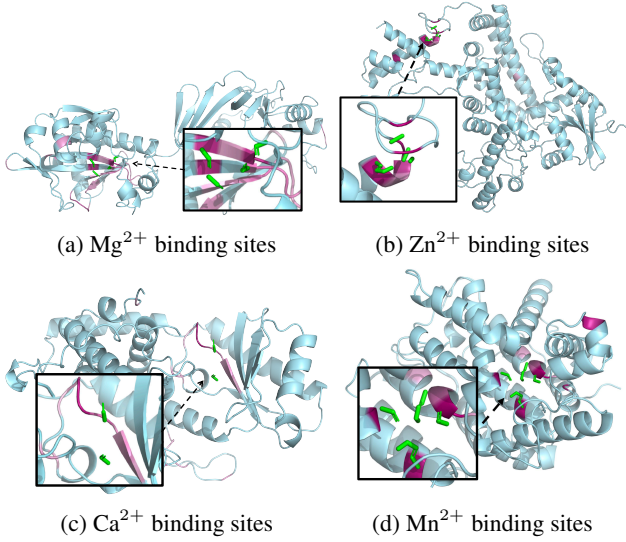


Figure 9: Interpretability results on metal ion binding-sites prediction

C.6 Relevance-based interpretation detailed results

We present detailed experimental results on three datasets in Figures 10, 11, and 12, analyzing the changes in AUC, Loss Recovered, Normalized MSE, and Proportion under different sparsity levels.

D Performance Analyses

D.1 Metrics

For evaluating the SAE performance, following (Rajamanoharan et al. 2024; Gao et al. 2024), we employ Nor-

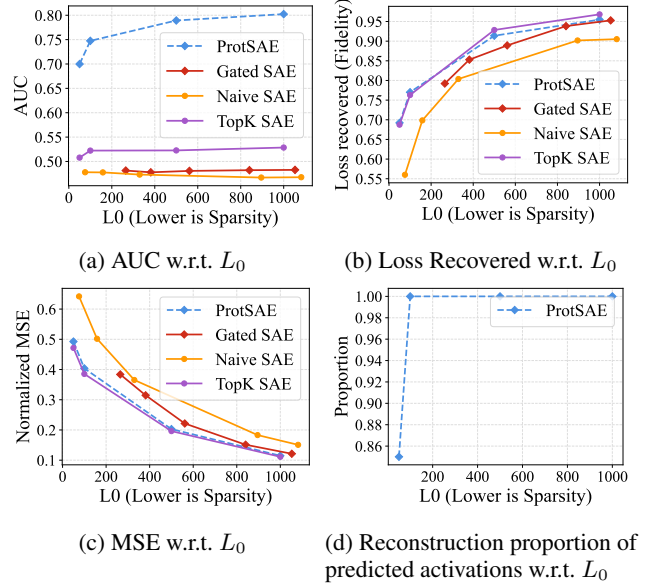


Figure 10: Performance comparison under different sparsity on the Molecular Function Ontology dataset

malized MSE and Loss Recovered to assess the reconstruction fidelity of the model, and L_0 to assess the sparsity. The formal definitions of these metrics are given below.

The L_0 of a SAE is defined by the average number of active features on a given input $\mathbb{E}_{x \sim \mathcal{D}} \|f(x)\|_0$.

The loss recovered of an SAE is calculated from the average cross-entropy loss of the language model on an evaluation dataset, when the SAE’s reconstructions are spliced into it. If we denote by $CE(\phi)$ the average loss of the language model when we splice in a function $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^n$ at the SAE’s site during the model’s forward pass, then loss recovered is

$$1 - \frac{CE(\hat{x} \circ f) - CE(\text{Id})}{CE(\zeta) - CE(\text{Id})}, \quad (41)$$

where $\hat{x} \circ f$ is the autoencoder function, $\zeta : \mathbf{x} \mapsto \mathbf{0}$ the zero-ablation function, and $\text{Id} : \mathbf{x} \mapsto \mathbf{x}$ the identity function.

We compute the mean squared error MSE between reconstructed output \hat{x} and the ground-truth x , normalized by the MSE obtained using the mean \bar{x} as reconstruction:

$$\text{Normalized_MSE}(\hat{x}, x) = \frac{\text{MSE}(\hat{x}, x)}{\text{MSE}(\bar{x}, x)}. \quad (42)$$

D.2 Detailed results

We present experimental results on three protein function prediction datasets in Figures 10, 11, and 12, analyzing the changes in AUC, Loss Recovered, Normalized MSE, and Proportion under different sparsity levels. ProtSAE consistently outperforms baselines in AUC under different sparsity levels, showing strong semantic retention. Meanwhile, it maintains competitive reconstruction quality, confirming its ability to capture meaningful concepts without losing essential concept information. Furthermore, the higher reconstruction proportion of predicted activations with respect to

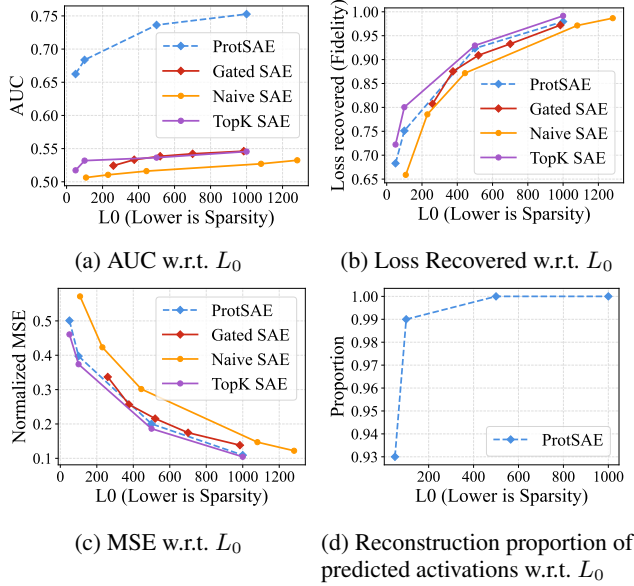


Figure 11: Performance comparison under different sparsity on the Biological Process Ontology dataset

L_0 indicates that the defined neurons effectively contribute to the reconstruction process.

E Steering Experiment

E.1 Settings

We design intervention experiments based on seven concepts to evaluate whether ProtSAE captures semantics related to specific biological concepts and consequently influences the generation of PLMs. The specific GO terms being intervened include: GO:0000793 “condensed chromosome”, GO:0009570 “chloroplast stroma”, GO:0036464 “cytoplasmic ribonucleoprotein granule”, GO:0009617 “response to bacterium”, GO:0016579 “protein deubiquitination”, GO:0032543 “mitochondrial translation”, and GO:0070646 “protein modification by small protein removal”. For each concept, we construct an evaluation dataset consisting of 5,000 samples, and generate 500 samples with intervention. Following previous works (Zongying et al. 2024; Lv et al. 2024; Liu et al. 2025), we assess the structural similarity between proteins generated after intervention and those in the evaluation dataset. A higher degree of structural similarity indicates that the protein language model, when guided by ProtSAE, tends to generate proteins structurally aligned with the targeted concept, potentially sharing similar functions. We use ESMFold (Lin et al. 2023b) for structure prediction of the generated protein sequences.

Metrics. We use Foldseek (van Kempen et al. 2022) to evaluate structural similarity using the Template Modeling score (TM-score) (Zhang and Skolnick 2004) and Root Mean Square Distance (RMSD) (Betancourt and Skolnick 2001). A higher TM-score and a lower RMSD, indicates greater structural similarity between the generated proteins

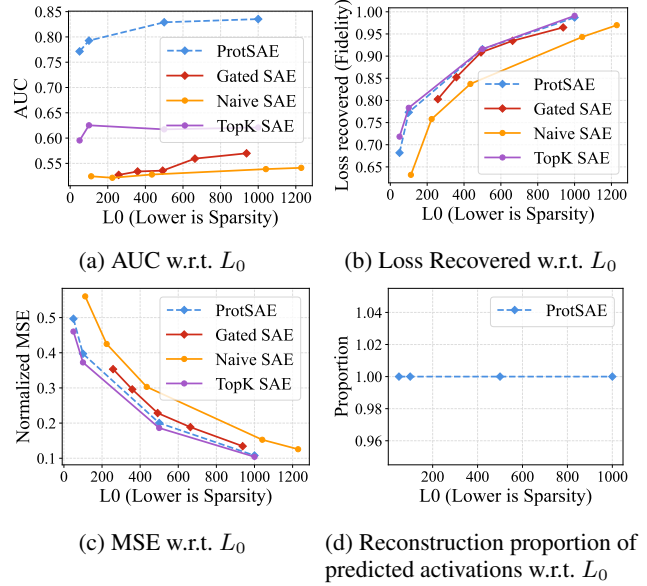


Figure 12: Performance comparison under different sparsity on the Cellular Component Ontology dataset

and those in the evaluation dataset. The predicted Local Distance Difference Test (pLDDT) score is used to evaluate the confidence of structure prediction. A higher pLDDT reflects more reliable predictions and greater structural stability.

Intervention method. Let $p = \{t_1, t_2, \dots, t_L\}$ be a protein sequence of length L , where t_j denotes the token at position j . Given a target concept c_i , we compute the relevance score $\pi_{\text{pred}}^{(i,j)}$ for each token t_j with respect to c_i , from Eq. (6). We define a masking threshold θ_i , as the 50%-position value in the sorted predicted importance scores. Then, we construct the masked sequence \tilde{p}_{c_i} by replacing tokens whose importance scores are below the threshold:

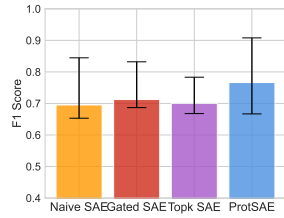
$$\theta_i = \text{Median} \left(\left\{ \pi_{\text{pred}}^{(i,j)} \right\}_{j=1}^L \right), \quad (43)$$

$$\tilde{p}_{c_i} = \left\{ t'_j \mid t'_j = \begin{cases} [\text{MASK}], & \text{if } \pi_{\text{pred}}^{(i,j)} < \theta_i \\ t_j, & \text{otherwise} \end{cases} \right\}. \quad (44)$$

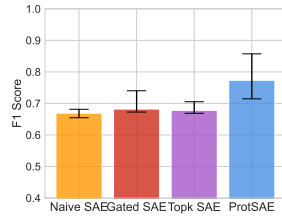
We use \tilde{p}_{c_i} into the protein language model (e.g., ESM2-15B) to obtain a reconstructed sequence \hat{p} . We simultaneously intervene on the target concept and its ancestor concepts. We observe that overly strong interventions may disrupt the semantic representations of PLMs. To mitigate this, for each concept, we enhance the activation by adding 1.2 times the mean of its Top- K activations.

E.2 Detailed Steering Results

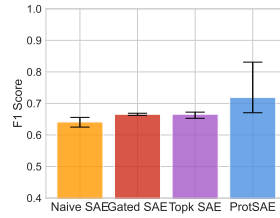
Here we show the average results before and after intervention on TM-score, RMSD and pLDDT in Figure 14. We observe that, after intervention, the generated samples exhibit significantly improved structural similarity to the evaluation set. At the same time, the structural stability of the generated



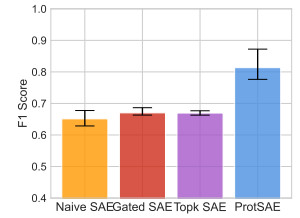
(a) Chloroplast stroma



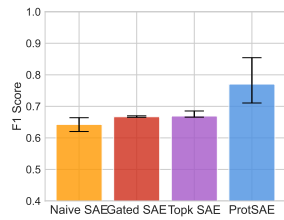
(b) Golgi membrane



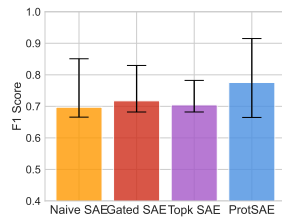
(c) Microtubule binding



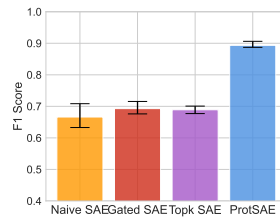
(d) Microtubule bundle formation



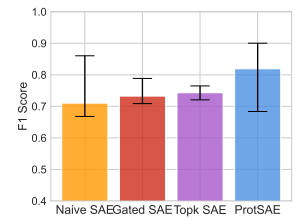
(e) Nuclear chromosome



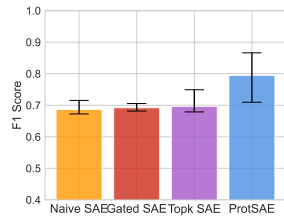
(f) Plastid envelope



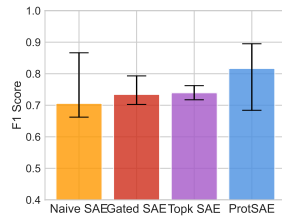
(g) Protein localization to cilium



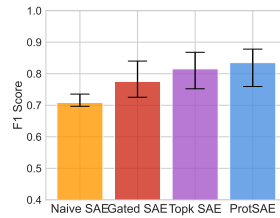
(h) Receptor ligand activity



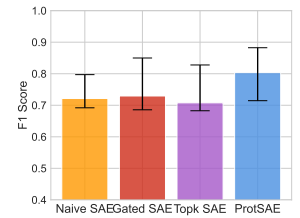
(i) rRNA metabolic process



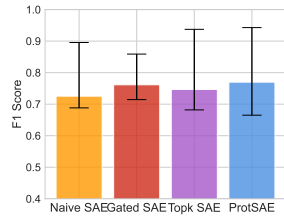
(j) Signaling receptor activator activity



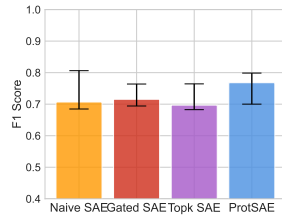
(k) Sodium ion transport



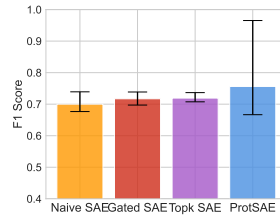
(l) Storage vacuole



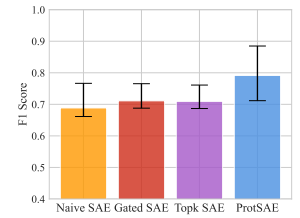
(m) Structural constituent of ribosome



(n) tRNA metabolic process



(o) UDP-glycosyltransferase activity



(p) Average performance on 15 concepts

Figure 13: Top-10 activated concept analysis results for 15 GO terms

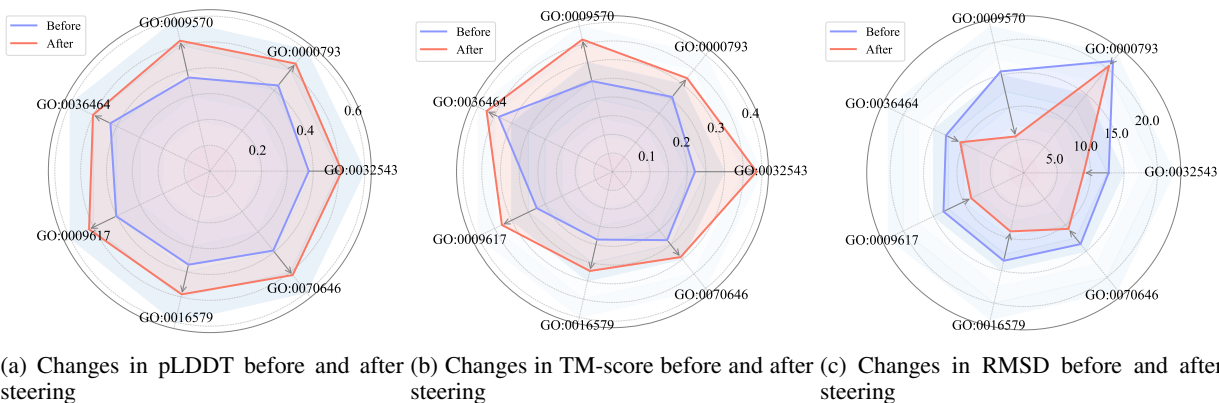


Figure 14: Detailed steering performance

Method	Dataset	L_0	Proportion	MSE	AUC	F_{\max}	S_{\min}	AUPR	Recovered
Naive SAE	MFO	329.6	1.00	0.365	0.473	0.332	14.686	0.239	0.804
Naive SAE	MFO	158.3	1.00	0.502	0.478	0.334	14.933	0.245	0.698
Naive SAE	MFO	75.7	1.00	0.642	0.478	0.335	15.159	0.246	0.560
Naive SAE	MFO	1081.7	1.00	0.151	0.468	0.327	14.578	0.225	0.905
Naive SAE	MFO	895.6	1.00	0.183	0.467	0.329	14.574	0.241	0.902
Gated SAE	MFO	1052.8	1.00	0.121	0.483	0.351	14.292	0.259	0.953
Gated SAE	MFO	840.6	1.00	0.151	0.482	0.400	14.314	0.254	0.938
Gated SAE	MFO	561.3	1.00	0.221	0.481	0.667	14.361	0.261	0.889
Gated SAE	MFO	381.6	1.00	0.315	0.478	0.340	14.439	0.233	0.853
Gated SAE	MFO	265.0	1.00	0.384	0.481	0.335	14.522	0.230	0.792
TopK SAE	MFO	50.0	1.00	0.471	0.508	0.500	14.177	0.264	0.687
TopK SAE	MFO	100.0	1.00	0.385	0.522	0.387	14.166	0.271	0.763
TopK SAE	MFO	500.0	1.00	0.196	0.523	0.364	14.140	0.268	0.928
TopK SAE	MFO	1000.0	1.00	0.111	0.529	0.355	14.131	0.277	0.968
ProtSAE	MFO	50.0	0.85	0.493	0.700	0.591	12.645	0.441	0.692
ProtSAE	MFO	100.0	1.00	0.404	0.747	0.591	12.352	0.434	0.770
ProtSAE	MFO	500.0	1.00	0.203	0.789	0.675	12.209	0.441	0.914
ProtSAE	MFO	1000.0	1.00	0.115	0.802	0.646	12.205	0.439	0.955

Table 4: Detailed results on the Molecular Function Ontology dataset

proteins also increases accordingly. This indicates that ProtSAE successfully captures protein structural features associated with the target function and can effectively guide the generation process through semantic intervention.

F Limitations

Currently, we evaluate ProtSAE on the ESM2 protein language model, and it has not been scaled to broader architectures or application scenarios. In future work, we aim to extend ProtSAE for more extensive exploration. For example, ProtSAE can be trained on the protein language models aligned with natural language. It enables us to investigate the association between natural language descriptions and structural fragments within protein sequences. ProtSAE can also be trained on the sequence-structure co-design model. Moreover, we can further explore how to extract interpretable features from protein structure prediction models and generative diffusion models.

G Tables of All Training Results

Here, we report the detailed training results of ProtSAE compared to SAE baselines across three datasets in Tables 4, 5, and 6.

Method	Dataset	L_0	Proportion	MSE	AUC	F_{\max}	S_{\min}	AUPR	Recovered
Naive SAE	BPO	443.8	1.00	0.302	0.516	0.306	43.743	0.236	0.872
Naive SAE	BPO	230.6	1.00	0.423	0.510	0.305	44.378	0.237	0.785
Naive SAE	BPO	108.1	1.00	0.571	0.506	0.305	44.223	0.237	0.659
Naive SAE	BPO	1281.5	1.00	0.122	0.532	0.305	43.566	0.228	0.987
Naive SAE	BPO	1080.4	1.00	0.147	0.527	0.305	43.570	0.224	0.971
Gated SAE	BPO	984.5	1.00	0.138	0.546	0.328	42.712	0.248	0.972
Gated SAE	BPO	700.4	1.00	0.175	0.542	0.326	42.775	0.249	0.933
Gated SAE	BPO	520.0	1.00	0.216	0.538	0.322	42.882	0.247	0.909
Gated SAE	BPO	375.8	1.00	0.257	0.533	0.318	43.040	0.243	0.875
Gated SAE	BPO	260.0	1.00	0.337	0.524	0.313	43.250	0.241	0.807
TopK SAE	BPO	50.0	1.00	0.461	0.517	0.324	42.850	0.257	0.722
TopK SAE	BPO	100.0	1.00	0.374	0.532	0.326	42.773	0.252	0.800
TopK SAE	BPO	500.0	1.00	0.186	0.536	0.324	42.792	0.249	0.929
TopK SAE	BPO	1000.0	1.00	0.104	0.546	0.325	42.800	0.246	0.992
ProtSAE	BPO	50.0	0.93	0.501	0.662	0.380	41.365	0.315	0.683
ProtSAE	BPO	100.0	0.99	0.397	0.684	0.381	41.249	0.316	0.751
ProtSAE	BPO	500.0	1.00	0.200	0.736	0.385	40.992	0.322	0.924
ProtSAE	BPO	1000.0	1.00	0.109	0.753	0.393	40.728	0.331	0.980

Table 5: Detailed results on the Biological Process Ontology dataset

Method	Dataset	L_0	Proportion	MSE	AUC	F_{\max}	S_{\min}	AUPR	Recovered
Naive SAE	CCO	434.6	1.0	0.303	0.528	0.631	11.824	0.572	0.837
Naive SAE	CCO	224.3	1.0	0.425	0.522	0.630	11.955	0.570	0.758
Naive SAE	CCO	111.7	1.0	0.560	0.524	0.625	12.000	0.600	0.632
Naive SAE	CCO	1228.8	1.0	0.126	0.541	0.631	11.744	0.571	0.970
Naive SAE	CCO	1040.9	1.0	0.153	0.539	0.630	11.793	0.571	0.943
Gated SAE	CCO	938.1	1.0	0.134	0.570	0.643	11.224	0.598	0.965
Gated SAE	CCO	663.6	1.0	0.189	0.559	0.641	11.292	0.596	0.934
Gated SAE	CCO	492.0	1.0	0.229	0.536	0.641	11.387	0.590	0.908
Gated SAE	CCO	358.7	1.0	0.296	0.534	0.638	11.470	0.584	0.852
Gated SAE	CCO	257.9	1.0	0.353	0.527	0.636	11.611	0.581	0.803
TopK SAE	CCO	50.0	1.0	0.460	0.595	0.651	11.170	0.608	0.718
TopK SAE	CCO	100.0	1.0	0.373	0.625	0.653	11.067	0.615	0.783
TopK SAE	CCO	500.0	1.0	0.186	0.617	0.651	11.131	0.611	0.916
TopK SAE	CCO	1000.0	1.0	0.104	0.621	0.651	11.110	0.614	0.991
ProtSAE	CCO	50.0	1.0	0.497	0.771	0.694	10.044	0.688	0.682
ProtSAE	CCO	100.0	1.0	0.397	0.793	0.696	9.978	0.692	0.773
ProtSAE	CCO	500.0	1.0	0.200	0.829	0.697	9.951	0.690	0.916
ProtSAE	CCO	1000.0	1.0	0.108	0.835	0.698	9.885	0.690	0.987

Table 6: Detailed results on the Cellular Component Ontology dataset