

Structure-Aligned Protein Language Model

Can (Sam) Chen^{1,2}, David Heurtel-Depeiges^{1,3,4}, Robert M. Vernon⁵, Christopher James Langmead⁵, Yoshua Bengio^{1,2}, and Quentin Fournier¹

¹Mila – Quebec AI Institute, ²Université de Montréal, ³Chandar Research Lab, ⁴Polytechnique Montréal, ⁵Amgen

Protein language models (pLMs) pre-trained on vast protein sequence databases excel at various downstream tasks but lack the structural knowledge essential for many biological applications. To address this, we integrate structural insights from pre-trained protein graph neural networks (pGNNs) into pLMs through a latent-level contrastive learning task. This task aligns residue representations from pLMs with those from pGNNs across multiple proteins, enriching pLMs with inter-protein structural knowledge. Additionally, we incorporate a physical-level task that infuses intra-protein structural knowledge by optimizing pLMs to predict structural tokens. The proposed *dual-task framework* effectively incorporates both inter-protein and intra-protein structural knowledge into pLMs. Given the variability in the quality of protein structures in PDB, we further introduce a *residue loss selection* module, which uses a small model trained on high-quality structures to select reliable yet challenging residue losses for the pLM to learn. Applying our structure alignment method to the state-of-the-art ESM2 and AMPLIFY results in notable performance gains across a wide range of tasks, including a 12.7% increase in ESM2 contact prediction. The data, code, and resulting SaESM2 and SaAMPLIFY models will be released on Hugging Face.

🤗 **Model Weights:** huggingface.co/chandar-lab/structure-alignment

🐙 **Code Repository:** github.com/chandar-lab/AMPLIFY

1. Introduction

Building on recent progress in natural language processing (Brown et al., 2020; Devlin et al., 2019), researchers have focused on pre-training protein language models (pLMs) on vast databases of protein sequences with masked language modeling (Rives et al., 2019; Hayes et al., 2024; Fournier et al., 2024) and next token prediction (Ferruz et al., 2022). These pLMs produce representations that demonstrate substantial potential across a variety of biological applications, including protein function annotation, enzyme-catalyzed reaction prediction, and protein classification (Hu et al., 2022).

While the sequence-only nature of pLMs contributes to their widespread adoption, they often struggle in tasks requiring detailed structural insights. For instance, ESM-GearNet outperforms ESM2 by 9.7% on the Human Protein-Protein Interaction classification task (Xu et al., 2022; Su et al., 2024). In this paper, we aim to develop a pLM that preserves its sequence-only nature for broader applicability yet is augmented with structural insights.

Given the availability of open-source pre-trained protein graph neural networks (pGNNs) (Zhang et al., 2023; Chen et al., 2023; Jumper et al., 2021), we investigate integrating pGNN-derived structural insights into pLMs. Specifically, we introduce a latent-level contrastive learning task for the structural alignment of pLMs. As illustrated in Figure 1, this task aligns residue hidden representations from the pLM (h_a) with those from the pGNN (h_g) across a batch of B proteins. During this process, the pGNN is frozen while the pLM is optimized to minimize the contrastive learning loss, enriching the pLM with inter-protein¹ structural knowledge. However, pure contrastive alignment may overemphasize residue-level patterns across the dataset, neglecting intra-protein¹ structural context (Zheng and Li, 2024). To address this, we introduce a physical-level task that trains the pLM to predict structural tokens z (representing physical conformations (van Kempen et al., 2022)) from its residue representations h_a . This reinforces the encoding of each residue within its protein, thereby enriching the pLM with intra-protein structural knowledge.

¹Note that “inter-protein” and “intra-protein” refer to tasks involving multiple proteins and within a single protein, respectively. This usage differs from the biological definition, where “inter-protein” refers to interactions between two proteins, and “intra-protein” refers to interactions within a single protein chain.

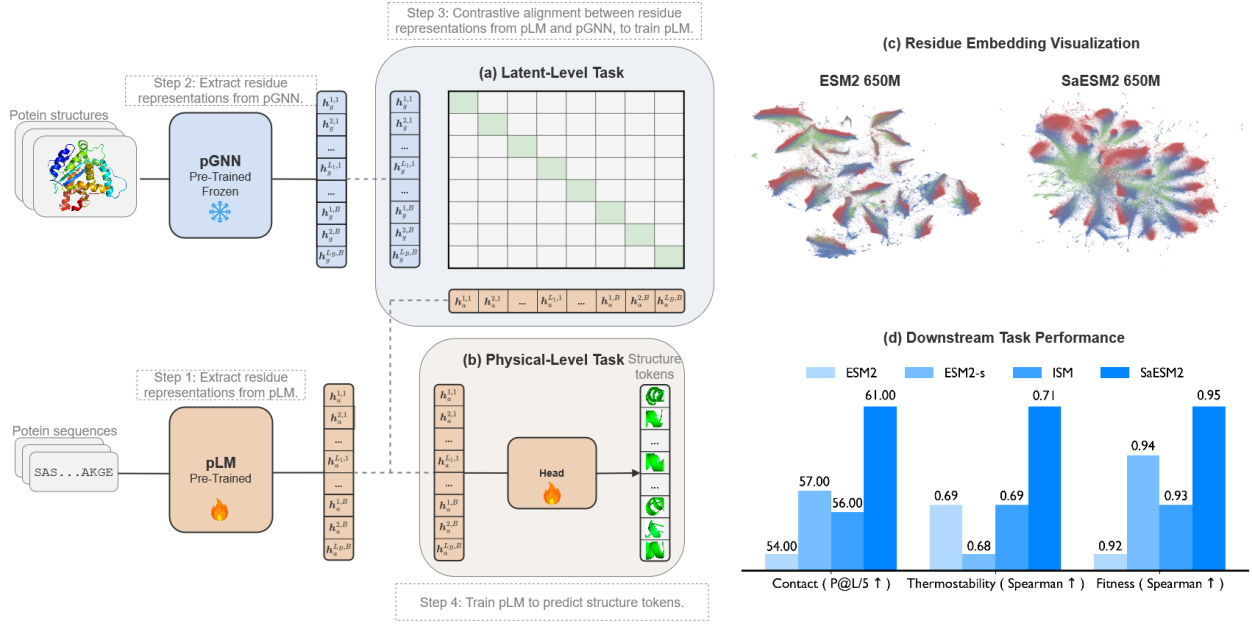


Figure 1: Overview of the *dual-task framework*. (a) Latent-level task: contrastively aligns residue representations from the pLM and pGNN, allowing the pLM to learn inter-protein structural knowledge. (b) Physical-level task: trains the pLM to predict structural tokens, incorporating intra-protein knowledge. (c) Residue embedding visualization: UMAP colored by secondary structure, showing that alignment improves separation. (d) Downstream task performance: structural knowledge improves contact map prediction, thermostability estimation, and fitness landscape modeling.

We combine latent and physical tasks, yielding three residue loss types for a batch of proteins with a total length N : (i) N sequence-to-structure contrastive losses from the latent-level task, (ii) N structure-to-sequence contrastive losses from the latent-level task, and (iii) N structure token prediction losses from the physical-level task. The *dual-task framework* effectively integrates inter-protein and intra-protein residue-level structural knowledge (§3.1). The masked language modeling loss is additionally incorporated to preserve the sequential knowledge of pLMs.

Given that some protein structure regions in the PDB are ambiguous or inaccurate (Burley et al., 2019), we propose a *residue loss selection* module that prioritizes residue losses aligned with high-quality protein structures across the $3 \times N$ total residue losses (§3.2). First, we use resolution and R-free metrics (Morris et al., 1992) to curate a high-quality reference set and train a small reference model on the set. Next, we compute the *excess loss*, defined as the difference between the residue loss of the current model and that of the reference model (Mindermann et al., 2022). Residue losses with high excess loss are selectively included in each loss type as they exhibit greater learnable potential. This module filters out inaccurate residues with high reference loss and easy residues with low current loss. By focusing on challenging yet reliable residue losses, the module improves both training effectiveness and efficiency.

We conducted 10 ablations to validate our design choices. Our analysis demonstrates that the proposed *dual-task framework* is crucial for optimal performance, with *residue loss selection* providing further gains. The models were evaluated on 6 tasks from xTrimoPGLM (Chen et al., 2024), 9 from SaProt (Su et al., 2023), and pseudo-perplexity using the high-quality held-out validation set from Fournier et al. (2024).

To summarize, our contributions are three-fold:

- We propose a *dual-task framework* that integrates inter-protein and intra-protein residue-level structural knowledge into pLMs.
- We develop a *residue loss selection* module that prioritizes challenging yet reliable residue losses, enhancing the learning process of pLMs.
- We conduct extensive experiments that demonstrate the effectiveness of our method.

2. Preliminaries

In this section, we introduce the preliminaries of protein language models, structure embeddings, and structure tokens used in this study, and provide a more detailed review of related work in [Appendix B](#).

2.1. Protein Language Models

Proteins can be represented as sequences of amino acids, where each amino acid a_i belongs to the set of the 20 common types. A protein sequence of length L is denoted as $\mathbf{a} = (a_1, a_2, \dots, a_L)$.

Protein language models are pre-trained on hundreds of millions of protein sequences using objectives such as masked language modeling (MLM) ([Hayes et al., 2024](#); [Rives et al., 2019](#)) and next-token prediction ([Ferruz et al., 2022](#)), capturing rich biophysical information. In this study, we focus on MLM-based pLMs, as proteins are not intrinsically left-to-right, and MLM has been shown to be highly effective for downstream tasks ([Lin et al., 2023b](#)).

A pre-trained pLM parameterized by θ is represented as $\text{pLM}(\cdot; \theta)$. The latent representation of a protein sequence \mathbf{a} is denoted as $\text{pLM}(\mathbf{a}; \theta) \in \mathbb{R}^{L \times D_a}$, where D_a is the embedding dimension. During pre-training, a subset of positions $\mathcal{M} \subset \{1, \dots, L\}$ is replaced with a [mask] token:

$$\tilde{a}_i = \begin{cases} [\text{mask}], & \text{if } i \in \mathcal{M}, \\ a_i, & \text{otherwise,} \end{cases} \quad (1)$$

where $\tilde{\mathbf{a}} = (\tilde{a}_1, \dots, \tilde{a}_L)$ represents the modified sequence. The model is trained to reconstruct the masked tokens by minimizing the masked language modeling loss:

$$\mathcal{L}_{\text{mlm}}(\theta, \alpha) = \frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} \ell_{\text{CE}}(\text{MLP}(\text{pLM}(\tilde{\mathbf{a}}; \theta)_i; \alpha), a_i), \quad (2)$$

where ℓ_{CE} is the cross-entropy loss, $\text{pLM}(\mathbf{a}; \theta)_i \in \mathbb{R}^{D_a}$ is the embedding at position i , and $\text{MLP}(\cdot; \alpha)$, parameterized by α , is a multi-layer perceptron head used during pre-training to predict the amino acid type.

2.2. Protein Structure Embeddings

Protein language models generate residue-level embeddings from protein sequences. In addition to the sequence perspective, proteins exist as 3D structures, and this physical nature largely determines their biological functions. Recent studies have also investigated deriving residue-level embeddings directly from protein 3D structures. One approach is to use the residue-level hidden representations generated by AlphaFold2 ([Jumper et al., 2021](#)), although their effectiveness for downstream tasks has since been questioned ([Hu et al., 2022](#)). GearNet ([Zhang et al., 2023](#)) addresses this limitation by pre-training a protein graph model encoder using multiview contrastive learning. Similarly, STEPS ([Chen et al., 2023](#)) improves protein structural representations by introducing multiple self-prediction tasks during graph model pre-training.

Given a protein graph \mathbf{g} , where each residue is a node and edges are defined based on both sequential and structural distances, and a pre-trained protein GNN model, the model outputs residue-level embeddings $\text{pGNN}(\mathbf{g}) \in \mathbb{R}^{L \times D_g}$, where L is the number of residues, and D_g is the embedding dimension of the graph-based residue representation.

2.3. Protein Structure Tokens

Inspired by the success of token-based protein language models, recent studies have explored the idea of tokenizing protein structures, representing a protein’s 3D conformation as a series of discrete structure tokens. A protein structure with L residues can be expressed as $\mathbf{z} = (z_1, z_2, \dots, z_L)$, where z_i denotes the structure token for the i -th residue.

Foldseek ([van Kempen et al., 2022](#)) introduces an efficient method for tokenizing protein structures, where each residue i is described by its geometric conformation relative to its spatially closest residue j . While this approach

has significantly accelerated homology detection, it incurs substantial information loss, thereby limiting its applicability for tasks requiring detailed structural reconstruction. To address this limitation, ProToken (Lin et al., 2023a) employs a symmetric encoder-decoder architecture that enables high-fidelity reconstruction of protein structures from tokens. Despite this advancement, these tokens have shown limited effectiveness in broader downstream applications Zhang et al. (2024a).

Recently, Hayes et al. (2024) developed an effective vector quantization variational autoencoder (VQ-VAE) tokenizer and integrated structure and sequence into a multi-modal large language model called ESM3. This approach effectively combines both modalities, improving the model’s versatility. We did not compare to ESM3 due to licensing restrictions. AIDO (Zhang et al., 2024a) further enhances structure tokenization by introducing a novel VQ-VAE with an equivariant encoder and an invariant decoder, ensuring a more robust representation of protein structures.

3. Method

3.1. Dual-Task Framework

We present our *dual-task framework*, consisting of a latent-level task and a physical-level task.

Latent-Level Task To incorporate structural insights from pre-trained pGNNs, we propose a latent-level contrastive learning task for the structure alignment of pLMs. Assuming a batch contains B proteins, with a total of $N = \sum_{b=1}^B L_b$ residues, we perform contrastive learning across all residues. We denote the pLM hidden representation of the i -th residue from the b_1 -th protein sequence \mathbf{a}_{b_1} as $\text{pLM}(\mathbf{a}_{b_1}; \boldsymbol{\theta})_i$, and the pGNN hidden representation of the j -th residue from the b_2 -th protein structure \mathbf{g}_{b_2} as $\text{pGNN}(\mathbf{g}_{b_2})_j$. Note that the parametrization of the pGNN is omitted for brevity, as the pGNN is frozen during training while only the pLM parameters $\boldsymbol{\theta}$ are optimized.

To align these embeddings, we introduce two linear layers, $\mathbf{W}_a \in \mathbb{R}^{D_a \times D}$ and $\mathbf{W}_g \in \mathbb{R}^{D_g \times D}$, and a learnable scalar s , parameterized as $\mathbf{W} = [\mathbf{W}_a; \mathbf{W}_g; s]$, which project both embeddings into the same dimension D . The similarity score between residues is computed as:

$$\delta(i, b_1, j, b_2) = s(\text{pLM}(\mathbf{a}_{b_1}; \boldsymbol{\theta})_i \mathbf{W}_a)^\top (\text{pGNN}(\mathbf{g}_{b_2})_j \mathbf{W}_g), \quad (3)$$

where s follows the approach in CLIP (Radford et al., 2021). In our experiments, we primarily use GearNet (Zhang et al., 2023) as the pGNN, pre-trained on the AlphaFold2 database (Varadi et al., 2022). We also evaluate the Evoformer within AlphaFold2 (Jumper et al., 2021) but find GearNet to be more effective for our purpose.

The sequence-to-structure residue contrastive loss for the i -th residue in the b_1 -th protein is:

$$\mathcal{L}_{a2g}(\boldsymbol{\theta}, \mathbf{W}, i, b_1) = -\log \frac{\exp(\delta(i, b_1, i, b_1))}{\sum_{b_2=1}^B \sum_{j=1}^{L_{b_2}} \exp(\delta(i, b_1, j, b_2))}. \quad (4)$$

The sequence-to-structure contrastive loss for the batch is then:

$$\mathcal{L}_{a2g}(\boldsymbol{\theta}, \mathbf{W}) = \frac{1}{N} \sum_{b_1=1}^B \sum_{i=1}^{L_{b_1}} \mathcal{L}_{a2g}(\boldsymbol{\theta}, \mathbf{W}, i, b_1). \quad (5)$$

A similar residue loss, $\mathcal{L}_{g2a}(\boldsymbol{\theta}, \mathbf{W}, b_2, j)$, can be defined for structure-to-sequence contrast, leading to the overall structure-to-sequence loss $\mathcal{L}_{g2a}(\boldsymbol{\theta}, \mathbf{W})$. The final latent-level loss is then given by:

$$\mathcal{L}_{\text{latent}}(\boldsymbol{\theta}, \mathbf{W}) = \frac{1}{2} (\mathcal{L}_{a2g}(\boldsymbol{\theta}, \mathbf{W}) + \mathcal{L}_{g2a}(\boldsymbol{\theta}, \mathbf{W})), \quad (6)$$

which enhances the pLM by integrating inter-protein residue-level structural knowledge.

Physical-Level Task However, pure contrastive alignment may overly emphasize residue-level structural patterns relative to the broader dataset, neglecting the intra-protein structural context. To address this, we introduce a physical-level task to reinforce the encoding of residue structure relative to its own protein.

This task trains the pLM to use the residue hidden representation to predict its structural token \mathbf{z} , which represents the residue’s physical conformation (van Kempen et al., 2022). The structure token prediction loss for the i -th residue in the b_1 -th protein is defined as:

$$\mathcal{L}_{\text{physical}}(\boldsymbol{\theta}, \boldsymbol{\beta}, i, b_1) = \ell_{\text{CE}}\left(\text{MLP}(\text{pLM}(\mathbf{a}_{b_1}; \boldsymbol{\theta})_i; \boldsymbol{\beta}), \mathbf{z}_{i, b_1}\right), \quad (7)$$

where ℓ_{CE} denotes the cross-entropy loss, MLP represents a multi-layer perceptron, and $\boldsymbol{\beta}$ are the parameters of the MLP. The overall physical-level loss is given by:

$$\mathcal{L}_{\text{physical}}(\boldsymbol{\theta}, \boldsymbol{\beta}) = \frac{1}{N} \sum_{b_1=1}^B \sum_{i=1}^{L_{b_1}} \mathcal{L}_{\text{physical}}(\boldsymbol{\theta}, \boldsymbol{\beta}, i, b_1), \quad (8)$$

infusing the pLM with intra-protein residue-level structural knowledge.

Overall Loss In addition to the dual-task losses, we incorporate the original MLM loss to preserve the sequential knowledge of pLMs, resulting in the final loss function:

$$\mathcal{L}_{\text{overall}}(\boldsymbol{\theta}, \boldsymbol{\alpha}, \mathbf{W}, \boldsymbol{\beta}) = \mathcal{L}_{\text{mlm}}(\boldsymbol{\theta}, \boldsymbol{\alpha}) + \gamma_{\text{latent}} \mathcal{L}_{\text{latent}}(\boldsymbol{\theta}, \mathbf{W}) + \gamma_{\text{physical}} \mathcal{L}_{\text{physical}}(\boldsymbol{\theta}, \boldsymbol{\beta}), \quad (9)$$

where γ_{latent} and γ_{physical} are weighting factors set to 0.5, ensuring equal importance for the latent-level and physical-level tasks. The weights are normalized such that $\gamma_{\text{latent}} + \gamma_{\text{physical}} = 1.0$, maintaining a balance between sequence and structure losses.

3.2. Residue Loss Selection

To address the challenge posed by ambiguous or inaccurate protein structures in the PDB (Burley et al., 2019), we propose a *residue loss selection* module. This module ensures both effectiveness and efficiency prioritizes residue losses that align with high-quality protein structures.

Reference Set We begin by curating a high-quality reference set using resolution and R-free metrics (Morris et al., 1992). Structures with resolution below 2.0Å and R-free lower than 0.20 are selected as a clean reference set. We then train a smaller language model, either ESM2 150M or AMPLIFY 120M, on the reference set with the same loss in Equation 9 and denote the optimized reference model parameters as $\boldsymbol{\theta}^r$, $\boldsymbol{\alpha}^r$, \mathbf{W}^r , $\boldsymbol{\beta}^r$. The resulting reference model is used to assess the residue loss of the alignment corpus.

Excess Loss For each residue loss discussed in subsection 3.1, we compute the *excess loss*, defined as the difference between the residue loss of the current model and that of the reference model:

$$\begin{aligned} \mathcal{L}_{\text{a2g}}(\Delta, i, b_1) &= \mathcal{L}_{\text{a2g}}(\boldsymbol{\theta}, \mathbf{W}, i, b_1) - \mathcal{L}_{\text{a2g}}(\boldsymbol{\theta}^r, \mathbf{W}^r, i, b_1), \\ \mathcal{L}_{\text{g2a}}(\Delta, j, b_2) &= \mathcal{L}_{\text{g2a}}(\boldsymbol{\theta}, \mathbf{W}, j, b_2) - \mathcal{L}_{\text{g2a}}(\boldsymbol{\theta}^r, \mathbf{W}^r, j, b_2), \\ \mathcal{L}_{\text{physical}}(\Delta, i, b_1) &= \mathcal{L}_{\text{physical}}(\boldsymbol{\theta}, \boldsymbol{\beta}, i, b_1) - \mathcal{L}_{\text{physical}}(\boldsymbol{\theta}^r, \boldsymbol{\beta}^r, i, b_1). \end{aligned} \quad (10)$$

where $\mathcal{L}_{\text{a2g}}(\Delta, i, b_1)$, $\mathcal{L}_{\text{g2a}}(\Delta, j, b_2)$, and $\mathcal{L}_{\text{physical}}(\Delta, i, b_1)$ represent the residue excess loss for sequence-to-structure, structure-to-sequence, and physical tasks, respectively.

Loss Selection Residue losses with high excess loss are prioritized for inclusion in the training as they exhibit greater learnable potential. This effectively filters out inaccurate residues, which typically have high reference

model loss, and excludes easy residues with low current model loss. We introduce a selection ratio ρ , selecting $N\rho$ residue losses for each type of loss. Taking \mathcal{L}_{a2g} as an example, we rewrite Equation 5 as:

$$\mathcal{L}_{a2g}(\boldsymbol{\theta}, \boldsymbol{w}) = \frac{1}{N\rho} \sum_{b_1=1}^B \sum_{i=1}^{L_{b_1}} \mathbb{1}(\mathcal{L}_{a2g}(\Delta, i, b_1), \rho) \mathcal{L}_{a2g}(\boldsymbol{\theta}, \boldsymbol{w}, i, b_1), \quad (11)$$

where $\mathbb{1}(\mathcal{L}_{a2g}(\Delta, i, b_1), \rho)$ equals 1 if $\mathcal{L}_{a2g}(\Delta, i, b_1)$ ranks in the top ρ of all $\mathcal{L}_{a2g}(\Delta, i, b_1)$ values, and 0 otherwise. This selection process is applied similarly for the other two types of losses. By focusing on challenging yet reliable residue losses, the *residue loss selection* module improves overall training effectiveness and efficiency.

4. Experiments

4.1. Structure Alignment Details

Our alignment dataset is sourced from the PDB database, comprising 129,732 proteins, including 116,713 proteins designated for training and 13,019 for validation/reference. The training protocol is adapted from the AMPLIFY stage-2 configuration (Fournier et al., 2024) with several modifications. We extend the pre-training with 20 epochs on our alignment dataset, with the learning rate linearly warming up from 0 to the peak rate over the first two epochs, followed by a cosine decay schedule for the subsequent 18 epochs. The peak learning rate for the language model is set at 1×10^{-4} , as per the AMPLIFY standard, while other modules, such as the structural linear classifier and the contrastive learning module, are set at 1×10^{-3} . The selection ratio ρ is set as 0.8. We employ the Zero Redundancy Optimizer (ZeRO) with DeepSpeed for efficiency and use 8 H100 GPUs. The effective batch size is 4,096 samples at a sequence length of 2,048, with longer proteins being randomly truncated.

4.2. Baseline Models

We evaluate the following sequence-only baseline pLMs: (1) **ESM2**: the standard ESM2 650M model (Lin et al., 2022); (2) **AMPLIFY**: the standard AMPLIFY 350M model (Fournier et al., 2024); (3) **ESM2-s**: a variant of ESM2 fine-tuned for fold classification (Zhang et al., 2024b); (4) **ISM**: a variant of ESM2 optimized for structure token prediction (Ouyang-Zhang et al., 2024). We denote our structure-aligned ESM2 and AMPLIFY models as **SaESM2** and **SaAMPLIFY**.

4.3. Structure Prediction

Tasks Structure-aligned models are expected to capture more nuanced insights of protein structures. We consider the following structure prediction tasks from xTrimoPGLM Chen et al. (2024): (1) **Contact**: two residues are considered in contact if their C_β atoms lie within 8Å (Rao et al., 2019). We evaluate this task using Top L/5 accuracy, considering residue pairs with a sequence separation greater than 6 and a sequence length cutoff of 512. (2) **Fold Classification (Fold)**: classify each protein sequence into one of 1,195 fold classes (Hou et al., 2018), with accuracy as the evaluation metric. (3) **Secondary Structure (SS)**: assign each residue to one of three secondary structure types (Rao et al., 2019), using accuracy as the evaluation metric.

To assess the quality of the learned representations, we freeze the backbone model and train a linear head for 20 epochs using a batch size of 128. We use a learning rate of 1×10^{-3} , with betas set to (0.9, 0.95) and a weight decay of 0.01 (Fournier et al., 2024). The linear head has a hidden size of 128, following (Chen et al., 2024). The linear head operates on residue embeddings for the token-level task (SS), on the mean-pooled residue embedding for the sequence-level task (Fold), and on pairwise residue embedding for the Contact task. We further visualize residue embeddings with secondary structure labels to assess structural alignment effectiveness in Appendix C.

Analysis As shown in Table 1, SaESM2 and SaAMPLIFY outperform their respective base models on all structure prediction tasks as well as existing alignment baselines on two out of three tasks, improving Contact

P@L/5 by 13% for ESM2 and 42% for AMPLIFY. ESM2-s, directly trained on fold classification, outperforms our alignment method on the corresponding task.

Table 1: Results on structure prediction.

	Contact	Fold	SS
	P@L/5 (↑)	Acc (↑)	Acc (↑)
ESM2	54.14	56.13	81.19
ESM2-s	57.38	67.26	79.87
ISM	55.55	47.87	83.24
SaESM2 (ours)	61.02	60.73	84.91
AMPLIFY	26.31	33.63	79.86
SaAMPLIFY (ours)	37.34	49.01	84.12

Table 2: Results on mutation effect prediction.

	Fluorescence	Fitness (GB1)	Stability
	Sp. (↑)	Sp. (↑)	Sp. (↑)
ESM2	0.687	0.916	0.742
ESM2-s	0.688	0.942	0.726
ISM	0.694	0.933	0.759
SaESM2	0.695	0.951	0.787
AMPLIFY	0.684	0.922	0.734
SaAMPLIFY	0.694	0.924	0.756

4.4. Mutation Effect Prediction

Tasks We evaluate our models on protein mutation effect prediction. Specifically, we consider three supervised tasks adopted in xTrimoPGLM (Chen et al., 2024): (1) **Fluorescence**: predicting the fluorescence intensity of green fluorescent protein mutants (Rao et al., 2019); (2) **Fitness (GB1)**: predicting the binding fitness of GB1 following mutations at four specific positions; (3) **Stability**: predicting relative protease resistance as a proxy measurement for stability. For this task, evaluation is performed on one-mutation neighborhoods of the most promising proteins (Rao et al., 2019). We report performance using the Spearman correlation coefficient (denoted as Sp.). The setup is the same as in subsection 4.3, except that we fine-tune the backbone with a learning rate of 1×10^{-4} .

Analysis As shown in Table 2, SaESM2 demonstrates a clear advantage over other models in mutation effect prediction, particularly in Fitness (GB1) and Stability. We also observe that AMPLIFY models achieve performance comparable to ESM2 despite being almost half the size, highlighting the importance of preserving the natural distribution of proteins in the pre-training corpus.

4.5. Property Prediction

Tasks We evaluate SaESM2 and SaAMPLIFY on a suite of downstream property prediction tasks (Xu et al., 2022; Dallago et al., 2021), which rely on structural information to some extent but are not direct structure prediction tasks. These include predictions of thermostability, metal ion binding, protein localization (DeepLoc), enzyme commission numbers (EC), gene ontology annotations (GO), and protein-protein interactions (HumanPPI). We follow the data splits and training protocols from SaProt (Su et al., 2024). While ESM2-based models are directly compatible with the SaProt codebase, AMPLIFY-based models are evaluated using a custom pipeline on the same splits, following the setup in subsection 4.4.

Analysis SaESM2 outperforms ESM2 across 6 out of 9 downstream property prediction tasks, underscoring the advantages of integrating structural alignment into pLMs beyond structure-only downstream applications. Furthermore, SaESM2 achieves state-of-the-art performance among sequence-only pLMs, including structure-aligned baselines, on these 6 tasks.

4.6. Pseudo Perplexity

To measure the impact of the structure alignment on the pLMs sequence-level knowledge, we compute the pseudo perplexity distributions of ESM2, AMPLIFY, and their structure-aligned variants as defined in Section 1.2.2 of Lin et al. (2022) using the validation set from Fournier et al. (2024). This set includes proteins with experimental evidence from reference proteomes based on high-quality genomes across all three domains of life and is designed to better reflect the natural protein distribution.

Figure 2 reveals that our structure alignment does increase pseudo perplexity, indicating a trade-off in which structural integration slightly compromises sequence modeling. However, despite this, both SaESM2 650M and

Table 3: Results on property prediction.

Method	Thermostability	HumanPPI	Metal Bind	EC	GO (Fmax (↑))			DeepLoc (Acc% (↑))	
	Spearman (↑)	Acc% (↑)	Acc% (↑)	Fmax (↑)	MF	BP	CC	Subcell.	Binary
ESM2	0.694	84.24	70.83	0.864	0.670	0.473	0.470	82.09	91.96
ESM2-s	0.683	82.61	71.13	0.861	0.673	0.479	0.458	82.85	92.28
ISM	0.695	81.52	69.94	0.872	0.666	0.471	0.497	82.63	92.28
SaESM2	0.705	85.33	72.33	0.861	0.669	0.488	0.461	83.39	92.68
AMPLIFY	0.646	75.00	69.02	0.680	0.566	0.389	0.479	76.56	90.06
SaAMPLIFY	0.675	71.11	70.83	0.773	0.586	0.412	0.512	77.25	91.22

SaAMPLIFY 350M maintain competitive pseudo perplexity scores, suggesting that structure alignment largely preserves the original pLMs sequence-level knowledge.

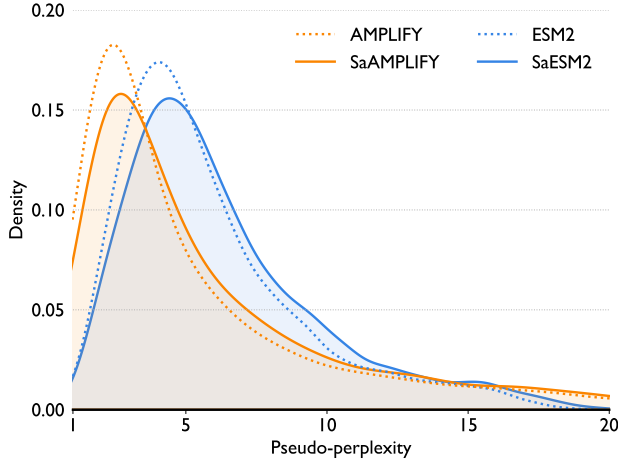


Figure 2: Pseudo perplexity distributions on the validation set introduced by Fournier et al. (2024).

4.7. Ablation Studies

We conduct extensive ablation studies on three tasks covering structure (Contact), mutation effect (Fluorescence), and property (Metal Bind) to evaluate the contribution of each design component.

Dual-Task Framework Our default setup employs a weighted combination of three losses: masked language modeling, latent-level, and physical-level, with weights [1, 0.5, 0.5], respectively. To assess the impact of each component, we experiment with the following configurations:

- *w/o latent*: Remove the latent-level loss, using weights [1, 0.0, 0.5].
- *w/o physical*: Remove the physical-level loss, using weights [1, 0.5, 0.0].
- *w/o dual*: Exclude both auxiliary losses, i.e. MLM fine-tuning on PDB, using weights [1, 0.0, 0.0].

Table 4: Ablations on dual-task framework.

	Contact	Fluorescence	Metal Bind
	P@L/5 (↑)	Spearman (↑)	Acc% (↑)
SaESM2 (all)	61.0	0.695	72.3
<i>w/o latent</i>	53.7 (−12.0%)	0.689 (−0.9%)	69.5 (−3.8%)
<i>w/o physical</i>	59.1 (−3.1%)	0.691 (−0.6%)	71.0 (−1.8%)
<i>w/o dual</i>	51.4 (−15.7%)	0.686 (−1.3%)	67.1 (−7.2%)
ESM2 (baseline)	54.1 (−11.3%)	0.687 (−1.2%)	70.8 (−2.1%)

As shown in Table 4, removing any loss term leads to performance degradation across all three tasks, confirming the effectiveness of our dual-task framework. Notably, the *w/o latent* setting performs worse than *w/o physical*, suggesting that the latent-level task contributes more significantly to the considered downstream tasks than the physical-level task. This supports our motivation that the physical-level task acts primarily as a structural constraint rather than a dominant learning signal.

Residue Loss Selection We compare our *residue-level selection* module with two alternative strategies that do not rely on reference models, instead selecting residues based solely on their individual loss values:

- *loss-large*: Select residues with high losses, assuming they offer greater learning potential.
- *loss-small*: Select residues with low losses, assuming they are cleaner and more accurate.

For comparison, we also include a *full* strategy that uses all residue losses without any selection.

Table 5: Ablations on residue loss selection.

	Contact	Fluorescence	Metal Bind
	P@L/5 (↑)	Spearman (↑)	Acc% (↑)
SaESM2 (<i>residue-loss selection</i>)	61.0	0.695	72.3
<i>loss-large</i>	60.6 (−0.7%)	0.693 (−0.3%)	71.3 (−1.4%)
<i>loss-small</i>	59.4 (−2.6%)	0.691 (−0.6%)	71.0 (−1.8%)
<i>full</i>	60.3 (−1.1%)	0.690 (−0.7%)	71.1 (−1.7%)
ESM2 (<i>baseline</i>)	54.1 (−11.3%)	0.687 (−1.2%)	70.8 (−2.1%)

As shown in Table 5, alternative selection strategies lead to decreased performance across all tasks, demonstrating the effectiveness of our *residue loss selection* module. While beneficial, its impact is less significant than that of the *dual-task framework*, likely due to the already high quality of the protein structures used and the extensive pre-training of base pLMs. We further visualize the validation loss curves for different loss selection strategies in §D, which further supports the superior effectiveness of our strategy.

Structure Embedding We further ablate the structure embeddings used in the latent-level task. In addition to our default GearNet embeddings (Zhang et al., 2023), we explore embeddings from the AlphaFold2 Evoformer model (Jumper et al., 2021), denoted as AF2. Specifically, we provide the protein structure as a template and perform only one Evoformer cycle to extract the embeddings to reduce computational cost.

Table 6: Ablations on structure embedding.

	Contact	Fluorescence	Metal Bind
	P@L/5 (↑)	Spearman (↑)	Acc% (↑)
SaESM2 (<i>GearNet</i>)	61.0	0.695	72.3
AF2	48.4 (−20.7%)	0.695 (−0.0%)	69.0 (−4.6%)
ESM2 (<i>baseline</i>)	54.1 (−11.3%)	0.687 (−1.2%)	70.8 (−2.1%)

As shown in Table 6, aligning to GearNet embeddings outperforms AF2 alignment on both Contact Prediction and Metal Bind tasks. Notably, AF2 alignment even underperforms the baseline ESM2 model without structural alignment. This observation is consistent with the findings of Hu et al. (2022), which suggest that embeddings from the AF2 are not well-suited for such downstream tasks.

Structure Token We further ablate the structure token used in the physical-level task. Our approach is based on *foldseek* structure tokens (van Kempen et al., 2022) and we explore *protoken* (Lin et al., 2023a) and *aido* (Zhang et al., 2024a), both of which employ a larger codebook size (512 compared to 20 for *foldseek*). We do not compare against the ESM3 structure token (Hayes et al., 2024) due to its strict commercial license.

Table 7: Ablations on structure token.

	Contact	Fluorescence	Metal Bind
	P@L/5 (↑)	Spearman (↑)	Acc% (↑)
SaESM2 (<i>foldseek</i>)	61.0	0.695	72.3
<i>protoken</i>	60.8 (−0.3%)	0.695 (+0.0%)	71.9 (−0.6%)
<i>aido</i>	61.9 (+1.5%)	0.695 (+0.0%)	70.5 (−2.5%)
ESM2	54.1 (−11.3%)	0.687 (−1.2%)	70.8 (−2.1%)

As shown in Table 7, *aido* outperforms *foldseek* on the contact prediction task, likely due to its finer-grained structural representation that injects richer structural insights into the pLM. *Protoken* performs slightly worse despite its larger codebook, likely due to *protoken* encoding global dependencies instead of emphasizing local neighborhoods like *foldseek*, which aligns more closely with our structure alignment approach. This observation is consistent with that of Zhang et al. (2024a). For the property prediction task Metal Bind, *foldseek* performs best, supporting the importance of local structure. All three tokens perform similarly on the fluorescence prediction task.

5. Conclusion

In this work, we propose to enrich sequence-only pLMs with structural knowledge. We incorporate structural insights from pre-trained pGNNs into pLMs via a latent-level task, aligning residue representations across models. To infuse intra-protein structural knowledge, we introduce a physical-level task that trains pLMs to predict structural tokens. Additionally, we propose a *residue loss selection* module that identifies and emphasizes challenging yet reliable residue losses to guide learning. We validate our structure alignment approach on two pLMs, ESM2 and AMPLIFY, demonstrating improved performance across diverse downstream tasks. These results suggest that structure alignment could become an indispensable component for future pLMs.

References

- Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne Longpre, Nathan Lambert, Xinyi Wang, Niklas Muenighoff, Bairu Hou, Liangming Pan, Haewon Jeong, et al. A survey on data selection for language models. *arXiv preprint arXiv:2402.16827*, 2024.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 2020.
- Stephen K Burley, Helen M Berman, Charmi Bhikadiya, Chunxiao Bi, Li Chen, Luigi Di Costanzo, Cole Christie, Ken Dalenberg, Jose M Duarte, Shuchismita Dutta, et al. Rcsb protein data bank: biological macromolecular structures enabling research and education in fundamental biology, biomedicine, biotechnology and energy. *Nucleic acids research*, 2019.
- Bo Chen, Xingyi Cheng, Pan Li, Yangli-ao Geng, Jing Gong, Shen Li, Zhilei Bei, Xu Tan, Boyan Wang, Xin Zeng, et al. xtrimopglm: unified 100b-scale pre-trained transformer for deciphering the language of protein. *Nature Method*, 2024.
- Can Chen, Jingbo Zhou, Fan Wang, Xue Liu, and Dejing Dou. Structure-aware protein self-supervised learning. *Bioinformatics*, 39, 2023.
- Christian Dallago, Jody Mou, Kadina E Johnston, Bruce J Wittmann, Nicholas Bhattacharya, Samuel Goldman, Ali Madani, and Kevin K Yang. Flip: Benchmark tasks in fitness landscape inference for proteins. *bioRxiv*, pages 2021–11, 2021.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter*

- of the association for computational linguistics: human language technologies, volume 1 (long and short papers), pages 4171–4186, 2019.
- Noelia Ferruz, Steffen Schmidt, and Birte Höcker. Protgpt2 is a deep unsupervised language model for protein design. *Nature communications*, 13(1):4348, 2022.
- Quentin Fournier, Robert M Vernon, Almer van der Sloot, Benjamin Schulz, Sarath Chandar, and Christopher James Langmead. Protein language models: Is scaling necessary? *bioRxiv*, 2024.
- Daria Frolova, Marina Pak, Anna Litvin, Ilya Sharov, Dmitry Ivankov, and Ivan Oseledets. Mulan: Multimodal protein language model for sequence and structure encoding. *bioRxiv*, pages 2024–05, 2024.
- Tomas Hayes, Roshan Rao, Halil Akin, Nicholas J Sofroniew, Deniz Oktay, Zeming Lin, Robert Verkuil, Vincent Q Tran, Jonathan Deaton, Marius Wiggert, et al. Simulating 500 million years of evolution with a language model. *bioRxiv*, pages 2024–07, 2024.
- Michael Heinzinger, Konstantin Weissenow, Joaquin Gomez Sanchez, Adrian Henkel, Milot Mirdita, Martin Steinegger, and Burkhard Rost. Bilingual language model for protein sequence and structure. *NAR Genomics and Bioinformatics*, 6(4):lqae150, 2024.
- Jie Hou, Badri Adhikari, and Jianlin Cheng. Deepsf: deep convolutional neural network for mapping protein sequences to folds. *Bioinformatics*, 2018.
- Mingyang Hu, Fajie Yuan, Kevin Yang, Fusong Ju, Jin Su, Hui Wang, Fei Yang, and Qiuyang Ding. Exploring evolution-aware &-free protein language models as protein function predictors. *Advances in Neural Information Processing Systems*, 2022.
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *nature*, 2021.
- Mingchen Li, Yang Tan, Xinzhu Ma, Bozitao Zhong, Huiqun Yu, Ziyi Zhou, Wanli Ouyang, Bingxin Zhou, Liang Hong, and Pan Tan. Prosst: Protein language modeling with quantized structure and disentangled attention. *bioRxiv*, pages 2024–04, 2024.
- Xiaohan Lin, Zhenyu Chen, Yanheng Li, Xingyu Lu, Chuanliu Fan, Ziqiang Cao, Shihao Feng, Yi Qin Gao, and Jun Zhang. Protokens: A machine-learned language for compact and informative encoding of protein 3d structures. *bioRxiv*. 2023a.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, et al. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv*, 2022.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 2023b.
- Zhenghao Lin, Zhibin Gou, Yeyun Gong, Xiao Liu, yelong shen, Ruochen Xu, Chen Lin, Yujiu Yang, Jian Jiao, Nan Duan, and Weizhu Chen. Not all tokens are what you need for pretraining. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=ONMzBwqaAJ>.
- Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- Soren Mindermann, Jan M Brauner, Muhammed T Razzak, Mrinank Sharma, Andreas Kirsch, Winnie Xu, Benedikt Höltgen, Aidan N Gomez, Adrien Morisot, Sebastian Farquhar, et al. Prioritized training on points that are learnable, worth learning, and not yet learnt. In *International Conference on Machine Learning*. PMLR, 2022.

- Anne Louise Morris, Malcolm W MacArthur, E Gail Hutchinson, and Janet M Thornton. Stereochemical quality of protein structure coordinates. *Proteins: Structure, Function, and Bioinformatics*, 1992.
- Jeffrey Ouyang-Zhang, Chengyue Gong, Yue Zhao, Philipp Krähenbühl, Adam R Klivans, and Daniel Jesus Diaz. Distilling structural representations into protein sequence models. *bioRxiv*, pages 2024–11, 2024.
- Daniel Penaherrera and David Ryan Koes. Structure-infused protein language models. *bioRxiv*, pages 2023–12, 2024.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter Abbeel, and Yun Song. Evaluating protein transfer learning with tape. *Advances in neural information processing systems*, 2019.
- Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *PNAS*, 2019. doi: 10.1101/622803. URL <https://www.biorxiv.org/content/10.1101/622803v4>.
- Louis Robinson, Timothy Atkinson, Liviu Copoiu, Patrick Bordes, Thomas Pierrot, and Thomas D Barrett. Contrasting sequence with structure: Pre-training graph representations with plms. *bioRxiv*, 2023.
- Jin Su, Chenchen Han, Yuyang Zhou, Junjie Shan, Xibin Zhou, and Fajie Yuan. Saprot: Protein language modeling with structure-aware vocabulary. *bioRxiv*, 2023.
- Jin Su, Chenchen Han, Yuyang Zhou, Junjie Shan, Xibin Zhou, and Fajie Yuan. Saprot: Protein language modeling with structure-aware vocabulary. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=6MRm3G4NiU>.
- Yuanfei Sun and Yang Shen. Structure-informed protein language models are robust predictors for variant effects. *Human Genetics*, 2024.
- Michel van Kempen, Stephanie S Kim, Charlotte Tumescheit, Milot Mirdita, Cameron LM Gilchrist, Johannes Söding, and Martin Steinegger. Foldseek: fast and accurate protein structure search. *Biorxiv*, pages 2022–02, 2022.
- Mihaly Varadi, Stephen Anyango, Mandar Deshpande, Sreenath Nair, Cindy Natassia, Galabina Yordanova, David Yuan, Oana Stroe, Gemma Wood, Agata Laydon, et al. Alphafold protein structure database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic acids research*, 2022.
- Minghao Xu, Zuobai Zhang, Jiarui Lu, Zhaocheng Zhu, Yangtian Zhang, Ma Chang, Runcheng Liu, and Jian Tang. Peer: a comprehensive and multi-task benchmark for protein sequence understanding. *Advances in Neural Information Processing Systems*, 35:35156–35173, 2022.
- Jiayou Zhang, Barthélemy Meynard-Piganeau, James Gong, Xingyi Cheng, Yingtao Luo, Hugo Ly, Le Song, and Eric Xing. Balancing locality and reconstruction in protein structure tokenizer. *bioRxiv*, pages 2024–12, 2024a.
- Zuobai Zhang, Minghao Xu, Arian Jamasb, Vijil Chenthamarakshan, Aurelie Lozano, Payel Das, and Jian Tang. Protein representation learning by geometric structure pretraining. In *International Conference on Learning Representations*, 2023.
- Zuobai Zhang, Jiarui Lu, Vijil Chenthamarakshan, Aurélie Lozano, Payel Das, and Jian Tang. Structure-informed protein language model. *arXiv preprint arXiv:2402.05856*, 2024b.
- Jiangbin Zheng and Stan Z Li. Ccpl: Cross-modal contrastive protein learning. In *International Conference on Pattern Recognition*, pages 22–38. Springer, 2024.

A. Terminology

B. Related Work

Structure Language Models There are two main types of structure language models. The first requires explicit structural input, such as structure tokens (Su et al., 2024; Heinzinger et al., 2024; Li et al., 2024) or torsion angles (Frolova et al., 2024). However, these models depend on potentially unreliable or inaccurate structural data, and protein structure databases like the PDB are much smaller than sequence-only databases. Additionally, many proteins lack a well-defined, rigid structure. The second type only requires protein sequences as input and integrates structural insights during pre-training. For example, Zhang et al. (2024b) introduces a physical-level task for fold prediction, though it is somewhat coarse. Sun and Shen (2024) proposes several physical-level tasks, including secondary structure and distance map predictions, to incorporate structural knowledge into the pLM, while Ouyang-Zhang et al. (2024) focuses on structure token prediction. Penaherrera and Koes (2024) uses a similar contrastive learning loss, but limits its focus to masked residues and does not utilize advanced pre-trained GNN models. Finally, AlphaFold2 (Jumper et al., 2021) and ESMFold (Lin et al., 2023b) use sequence encoders, namely Evoformer and ESM2, followed by structure prediction modules. However, their focus is on structure prediction, and AlphaFold2 embeddings have been shown to be less effective than ESM2 embeddings for downstream tasks (Hu et al., 2022).

While recent studies have explored how to incorporate knowledge from pre-trained pLMs into pGNNs (Zheng and Li, 2024; Chen et al., 2023; Robinson et al., 2023), their focus is on improving pGNNs rather than pLMs, and no prior work has explored integrating structural insights from pre-trained pGNNs into pLMs. Our work bridges this gap by introducing the latent-level task, thereby enriching the pLMs with comprehensive structural insights from pre-trained pGNNs.

Data Selection Data selection is a critical component in training protein models. AlphaFold2 (Jumper et al., 2021) filters proteins with a resolution higher than 9Å and excludes sequences where a single amino acid accounts for over 80% of the input sequence. Additionally, it samples protein chains based on length to rebalance distribution and cluster size to reduce redundancy, which risks deviating from the natural distribution shaped by evolutionary selection. ESM2 (Lin et al., 2023b) adopts comparable sampling strategies while AMPLIFY (Fournier et al., 2024) curates a validation set of proteins with experimental evidence at the protein or transcript level from reference proteomes derived from high-quality genomes across all three phylogenetic domains, aiming to better represent the natural protein distribution.

Data selection has also been extensively explored in natural language model pre-training, incorporating techniques such as filtering, heuristics, and domain-specific selection (Albalak et al., 2024). Our *residue loss selection* module is inspired by prior work (Lin et al., 2024), which uses excess loss to identify useful tokens in language pre-training. However, our approach differs significantly by operating at a finer granularity through residue-level loss. Given the multi-loss structure of our framework, where each residue incurs three types of losses, we focus on those with high excess loss in each specific category. Crucially, our work is rooted in the protein research rather than natural language, reflecting the unique challenges and requirements of protein modeling.

C. Residue Embedding Visualization

In order to qualitatively assess the effectiveness of our structure alignment technique, we visualize the residue embeddings extracted from the final layer of ESM2 and AMPLIFY before and after aligning them. Specifically, we analyze 1,000 proteins from the Secondary Structure task, where each residue is color-coded based on its annotation to one of three secondary structure labels. We use UMAP (McInnes et al., 2018) to project high-dimensional data into a two-dimensional space with 50 nearest neighbors.

As shown in Figure 3, applying structure alignment improves the discrimination between secondary structures. In particular, the aligned embeddings (SaESM2 and SaAMPLIFY) exhibit clearer separation compared to their unaligned counterparts. Additionally, Figure 4 shows that amino acids sharing similar physical properties are located closer in the embedding space for aligned models compared to the unaligned baseline.

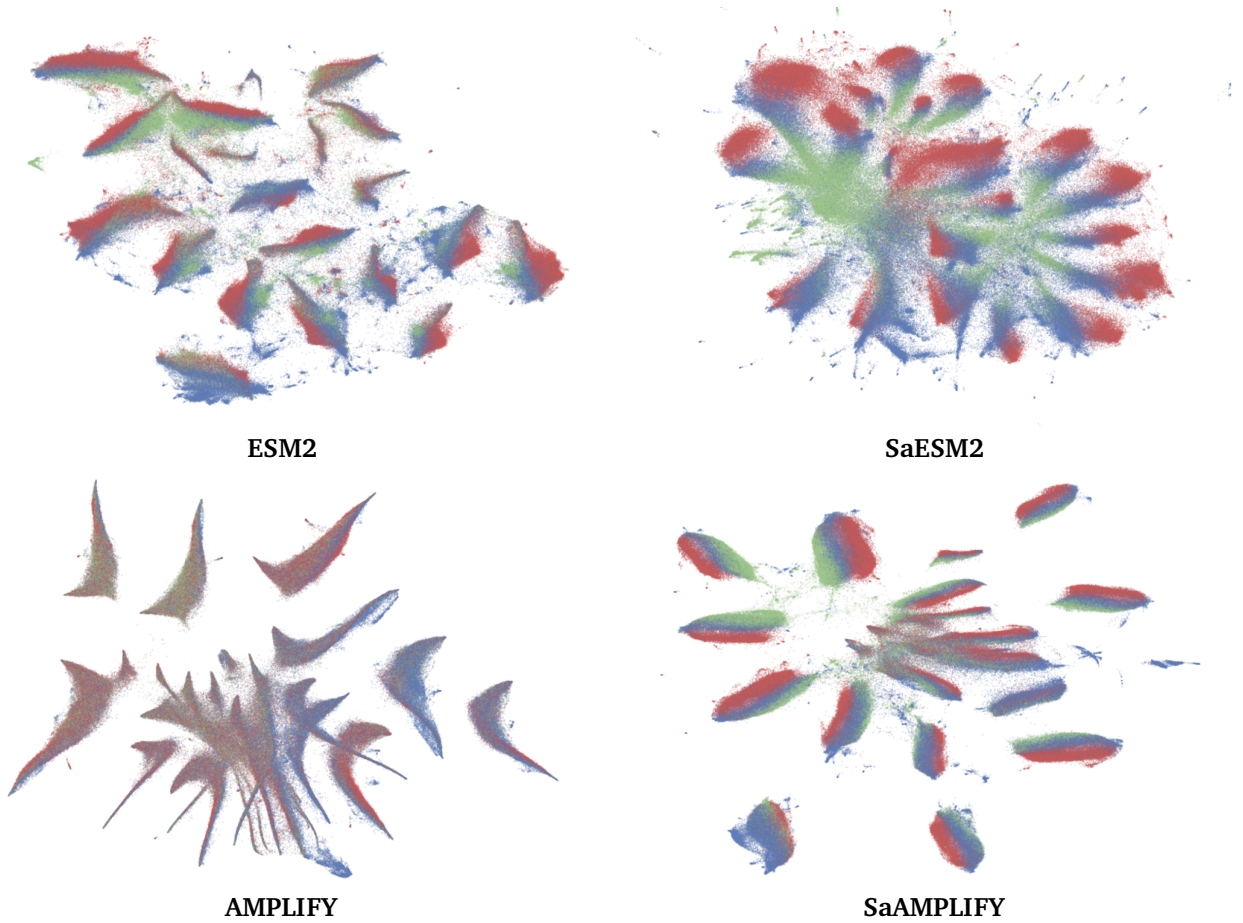


Figure 3: Residue embeddings colored by secondary structure type colored in blue, red, and green across four models: ESM2, SaESM2, AMPLIFY and SaAMPLIFY.

D. Loss Curve Analysis of Residue Loss Selection

To assess the effectiveness of our proposed *residue loss selection* module, we analyze validation loss curves across four strategies: ours, loss large, loss small, and full. These are shown in [Figure 5](#) (overall loss), [Figure 6](#) (MLM loss), [Figure 7](#) (latent-level loss), and [Figure 8](#) (physical-level loss). Recall that the overall loss is defined as:

$$\mathcal{L}_{\text{overall}} = \mathcal{L}_{\text{mlm}} + 0.5\mathcal{L}_{\text{latent}} + 0.5\mathcal{L}_{\text{physical}}. \quad (12)$$

As seen in [Figure 5](#), our strategy consistently achieves the lowest overall loss, demonstrating superior training effectiveness and efficiency. [Figure 7](#) shows that the primary reduction comes from the latent-level loss, indicating that our method successfully identifies informative and challenging latent-level residue losses to enhance learning. In contrast, [Figure 8](#) shows negligible differences in physical-level loss across most strategies, except for loss small. We attribute this to the limited Foldseek codebook size (20), which provides only coarse structural information, reducing the potential benefit of residue loss selection at this level. Notably, the loss small strategy results in high physical-level loss, likely due to its focus on easy-to-learn residues, which fail to contribute meaningful structural insights to the pLMs. We further experiment with joint training on both the training and validation sets. However, this led to degraded downstream performance, likely due to overfitting on the validation set.

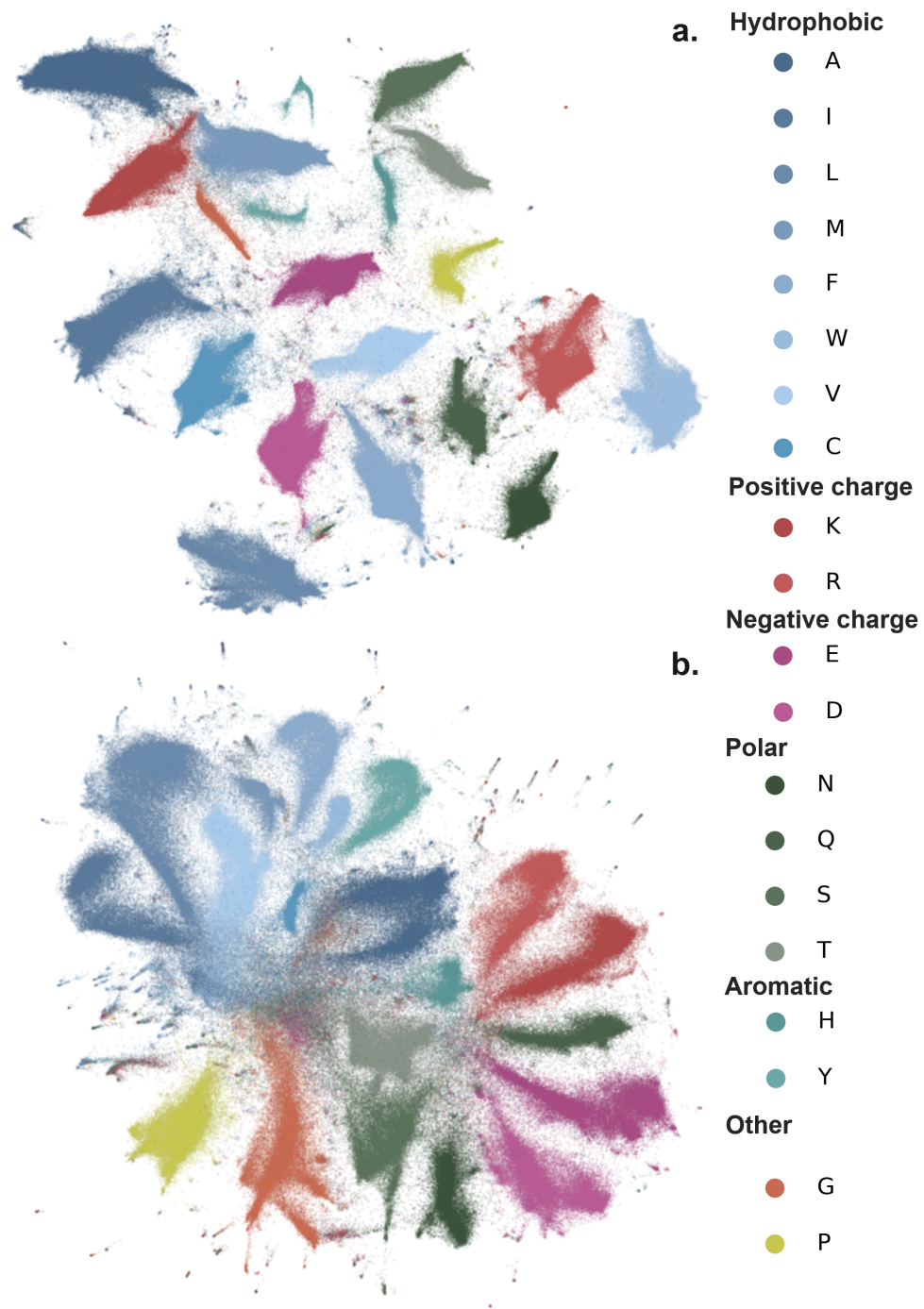


Figure 4: Residue embeddings colored by amino acid type for ESM2 (a.) and SaESM2 (b.). Amino acids with similar physical properties are colored with a gradient of the same color. Embeddings for SaESM2 clearly show a latent space more physically coherent.

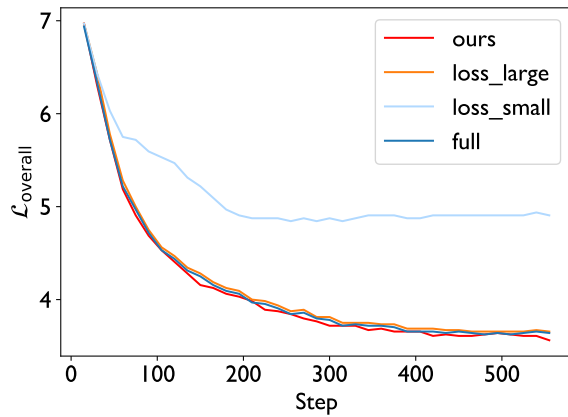


Figure 5: Overall loss.

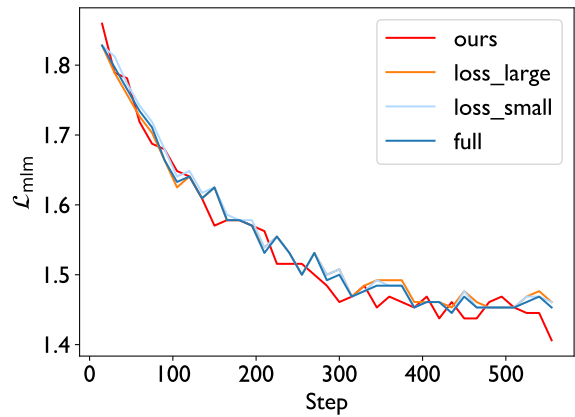


Figure 6: MLM loss.

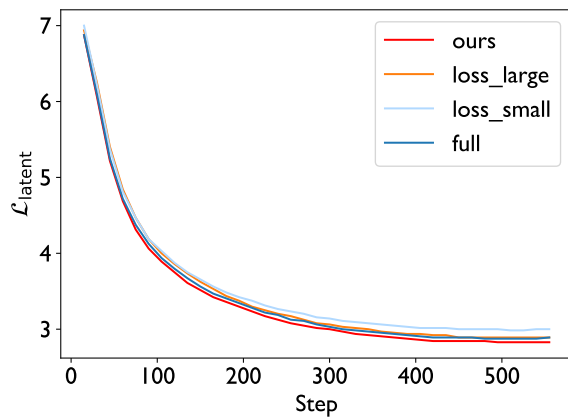


Figure 7: Latent-level loss.

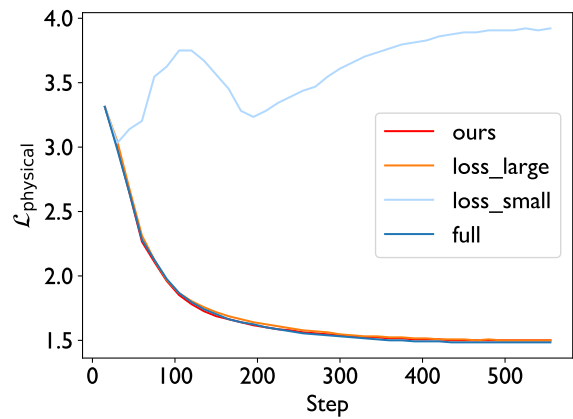


Figure 8: Physical-level loss.