

# PROTEIN SEQUENCE DOMAIN ANNOTATION USING LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Protein function inference relies on annotating protein domains via sequence similarity, often modeled through profile Hidden Markov Models (profile HMMs), which capture evolutionary diversity within related domains. However, profile HMMs make strong simplifying independence assumptions when modeling residues in a sequence. Here, we introduce PSALM (Protein Sequence Annotation using Language Models), a hierarchical approach that relaxes these assumptions and uses representations of protein sequences learned by protein language models to enable high-sensitivity, high-specificity residue-level protein sequence annotation. We also develop the Multi-Domain Protein Homology Benchmark (MDPH-Bench), a benchmark for protein sequence domain annotation, where training and test sequences have been rigorously split to share no similarity between any of their domains at a given threshold of sequence identity. Prior benchmarks, which split one domain family at a time, do not support methods for annotating multi-domain proteins, where training and test sequences need to have multiple domains from different families. We validate PSALM’s performance on MDPH-Bench and highlight PSALM as a promising alternative to HMMER, a state-of-the-art profile HMM-based method, for protein sequence annotation.

## 1 INTRODUCTION

Proteins are composed of distinct structural and functional units conserved through evolution, known as domains. The primary aim of protein sequence annotation is to locate and characterize these domains within a given sequence. Insight into the individual functions of these domains, which may act independently or in concert with neighboring domains, may shed light on the overall biological role of the protein (Fig. 1). Since experimental characterization of protein function can be difficult, function is often inferred and annotated through sequence similarity (homology) to domains with known function (Pearson, 2013; Eddy, 1998). As the size of protein sequence databases and the number of protein sequences with unknown function continue to grow rapidly (UniProt Consortium, 2023), methods for annotating protein sequences have been critical for exploiting this wealth of information about the molecular basis and evolutionary trajectory of life.

The state of the art in protein domain sequence annotation uses profile hidden Markov models (profile HMMs) to detect domains (Eddy, 2011) and profile/profile comparison to identify homologous domains (Remmert et al., 2012). Databases of protein domain families, like Pfam (Mistry et al., 2021), categorize millions of protein sequences into approximately 20,000 domains. The state of the art uses profile hidden Markov models (profile HMMs) to identify domains (Eddy, 2011) and profile/profile comparison to identify homologous domains (Remmert et al., 2012). Annotation can be achieved by using databases of protein families to compare a protein sequence against 20,000 profiles rather than scanning against millions of individual sequences. Domain-based annotation goes beyond classifying the whole protein sequence; it identifies the domain composition as well as the boundaries of each domain. This “domain annotation” requires annotation at the residue level, labeling each amino acid symbol in the sequence. Domain annotation helps with function inference and avoids the “transitive identification catastrophe”, where sequence-level annotations inferred from the presence of one domain can erroneously transfer between sequences due to homology of an unrelated domain (Doerks et al., 1998).

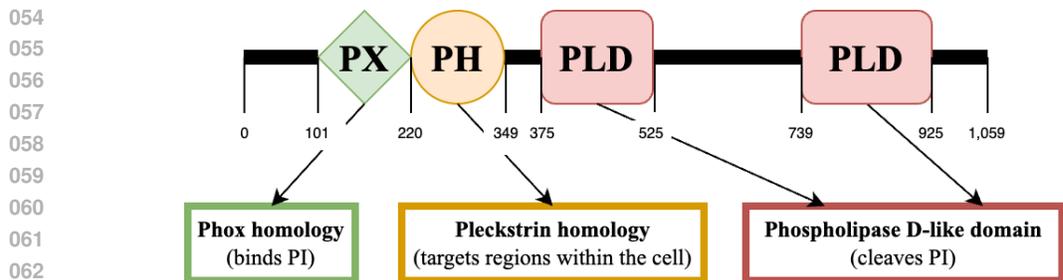


Figure 1: **Annotated domain architecture of a human phospholipase D1 protein (Q59EA4)** (EMBL-EBI, 2024; Paysan-Lafosse et al., 2023), featuring PX (phox), PH (pleckstrin homology), and PLD (phospholipase D-like) domains. Together, the function of these domains suggest that the full length protein (1,059 amino acids) is involved in phosphatidylcholine (PI) cleavage and intracellular signaling, consistent with experimental evidence.

Profile HMMs make simplifying independence assumptions when modeling residues in protein sequences, and there is considerable interest in developing more powerful methods that can better recognize distant evolutionary relationships with greater sensitivity. While most deep learning approaches to protein sequence similarity recognition focus on whole-protein or single-domain classification (Bileschi et al., 2022; Heinzinger et al., 2022; Nallapareddy et al., 2023; Kaminski et al., 2023; Hamamsy et al., 2023), they do not address the challenge of identifying individual domain subsequences within longer target sequences, which requires careful benchmarking and data curation to assess performance.

In this work, we introduce Protein Sequence Annotation with Language Models (PSALM), a novel approach that extends the capabilities of ESM-2, a pre-trained protein language model (pLM) (Lin et al., 2023), to predict *residue-level* sequence annotations. Our contributions include:

- **First deep learning model for residue-level protein domain annotation:** PSALM is the first deep learning approach to annotate domain boundaries and subsequences within multidomain proteins.
- **Relaxation of HMM independence assumptions for improved sensitivity and specificity:** PSALM leverages pLMs to overcome the simplifying assumptions of profile HMMs, allowing for greater sensitivity in detecting conserved domains across distantly-related sequences and higher specificity in identifying previously unannotated domains.
- **First benchmark for multidomain protein annotation, MDPH-Bench:** To enable robust evaluation, we introduce the Multi-Domain Protein Homology Benchmark (MDPH-Bench). This benchmark rigorously curates training and test sets to prevent any domain similarity above a predefined threshold, enabling realistic assessments of model performance across diverse domain families and multidomain proteins, which previous benchmarks do not support.

## 2 RELATED WORK

### 2.1 PROFILE HMMs

Profile HMMs use curated multiple sequence alignments (MSAs) of related domains, which reveal patterns of conservation and variability at the residue level, to model consensus using “match”, “insert”, and “delete” hidden states (Durbin et al., 1998; Eddy, 1998). Profile HMMs assume that the observed residues are conditionally independent given the hidden state. While this assumption simplifies the modeling process, it may limit the ability of profile HMMs to capture complex dependencies between residues in a sequence. These models serve as templates for comparison against the sequence of interest, enabling the identification of domains by finding subsequences that match the profile HMMs. Sequences with multiple, unrelated domains will require the use of multiple profile HMMs for annotation. HMMER (Eddy, 2011) is the state-of-the-art protein sequence domain annotation method and underlies many different databases, which organize related domains

108 into MSAs and profile HMMs at varying levels of granularity, enabling profile-based annotations at  
 109 the superfamily (Pandurangan et al., 2019), family (Mistry et al., 2021), and sub-family (Thomas  
 110 et al., 2022) levels.

## 111 2.2 DEEP MODELS

112 Recent efforts to apply deep learning methods to predict protein function from sequence have either  
 113 focused on predicting ontology-based functional annotation at the sequence level (Cao & Shen, 2021;  
 114 Hong et al., 2020; Kulmanov & Hoehndorf, 2020) or recognizing homology at the sequence level  
 115 (Heinzinger et al., 2022; Nallapareddy et al., 2023; Kaminski et al., 2023; Hamamsy et al., 2023). To  
 116 our knowledge, ProtENN, an ensemble of convolutional neural networks, represents the first attempt  
 117 to predict Pfam domains directly from protein sequences (Bileschi et al., 2022). ProtENN, however,  
 118 is constrained to make one domain prediction per input sequence and cannot natively identify domain  
 119 boundaries or multiple domains within a sequence without ad hoc post-processing. Additionally,  
 120 ProtENN cannot provide information on the contribution of an individual residue to a predicted  
 121 annotation.

## 122 3 METHODS

### 123 3.1 PROBLEM FORMULATION

124 Here, we formalize the residue-level sequence annotation problem as a mapping from a protein  
 125 sequence  $\mathbf{x} = (x_1, x_2, \dots, x_L)$  to a sequence of protein domain families  $\mathbf{y} = (y_1, y_2, \dots, y_L)$ . For  
 126 residue  $i$  in a sequence,  $x_i$  is an index  $1 \dots 25$  representing the  $i$ -th amino acid character (20 canonical,  
 127 2 non-canonical, and 3 ambiguous amino acid characters), and  $y_i$  is an index  $1 \dots D$  representing  
 128 the  $i$ -th protein domain family annotation, with  $D + 1$  for none. Approximately 23% of protein  
 129 sequences and 47% of residues across all sequences in UniProt do not belong to any Pfam domain  
 130 family (Mistry et al., 2021). The goal of residue-level annotation is to learn a model that predicts  
 131 domain family annotations for each residue in a protein sequence:

$$132 \hat{y}_i = \arg \max_f P(Y_i = f | \mathbf{x}), \quad (1)$$

133 where  $P(Y_i)$  is the distribution over  $D + 1$  family annotations for a given residue  $i$ .

### 134 3.2 PROTEIN LANGUAGE MODEL

135 Numerically encoding protein sequences is necessary as an input for machine learning tasks such  
 136 as classification. Protein language models (pLMs) learn vector representations of both individual  
 137 residues and full-length protein sequences, which capture long-range interactions, predict function  
 138 via transfer learning, and achieve state-of-the-art performance in several structure prediction tasks  
 139 (Bepler & Berger (2021); Rao et al. (2020); Meier et al. (2021); Elnaggar et al. (2021)). We use ESM-2  
 140 (specifically, the 8M, 35M, 150M, and 650M parameter models), a pre-trained general-purpose pLM  
 141 (Lin et al., 2023), to generate residue-level sequence embeddings  $\mathbf{x}'$  for a given sequence  $\mathbf{x}$ . ESM-2  
 142 was trained using a BERT-style masked token prediction task (Devlin et al., 2018), enabling it to  
 143 capture contextual information and dependencies within protein sequences and allowing us to replace  
 144  $\mathbf{x}$  with  $\mathbf{x}'$ .

### 145 3.3 PSALM

146 We introduce PSALM (Protein Sequence Annotation using Language Models), a method to predict  
 147 domains across a protein sequence at the residue-level. PSALM uses a hierarchical approach that  
 148 considers both individual protein domain families and clans, which are collections of evolutionarily  
 149 related (homologous) protein domain families categorized by Pfam (Finn et al., 2006). In Pfam 35.0,  
 150 approximately 45% of the 19,632 Pfam families are grouped into 655 clans, and a family can only  
 151 belong to at most one clan. While our primary aim is to predict protein domain families at each  
 152 residue, modeling clans – the super class – is an interpretable intermediate step that aids in identifying  
 153 areas of functional or structural importance that may not have clear family-level annotations.

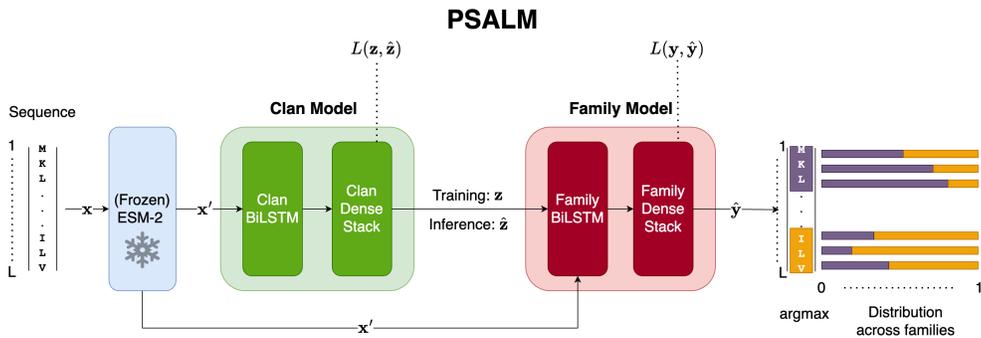


Figure 2: **Overview of residue-level protein sequence annotation with PSALM.** A sequence  $\mathbf{x}$  of length  $L$  is embedded as  $\mathbf{x}'$  with a frozen ESM-2. The PSALM clan and family models predict the clan annotations  $\hat{\mathbf{z}}$  and family annotations  $\hat{\mathbf{y}}$ , respectively, and are trained to minimize cross-entropy loss  $L(\cdot)$ . Here, the example outputs are predicted across a set of 2 families.

This intermediate annotation problem is a mapping from  $\mathbf{x}$  to a sequence of Pfam clans  $\mathbf{z} = (z_1, z_2, \dots, z_L)$ , where  $z_i$  is an index  $1 \dots C$  representing the  $i$ -th clan annotation, with  $C + 1$  for “non-clan” or  $C + 2$  for none. The non-clan annotation describes a residue which belongs to a domain family that is not a member of clan, and none refers to a residue which does not belong to a clan. For a given residue, the PSALM clan and family models learn to predict:

$$\hat{z}_i = \arg \max_c P(Z_i = c | \mathbf{x}') \quad (2)$$

$$\hat{y}_i = \arg \max_f P(Y_i = f | Z_i = C(f), \mathbf{x}') P(Z_i = C(f) | \mathbf{x}'), \quad (3)$$

where  $C(f)$  is the clan label to which family  $f$  belongs, and  $P(Z_i)$  is the distribution over all  $C + 2$  clan annotations for a given residue  $i$ . The inclusion of a separate clan prediction task ensures the interpretability of the clan model, preventing it from becoming an abstract hidden state. The family model is trained via the teacher forcing algorithm (Williams & Zipser, 1989), where it is provided the correct clan annotation for each residue in order to mitigate error propagation.

The clan and family models follow a similar structure and are trained separately. Protein sequences are initially embedded at the residue level using a pre-trained and frozen instance of ESM-2. The resulting embeddings are then passed into a bidirectional Long Short-Term Memory (BiLSTM) layer to capture sequential dependencies (Hochreiter & Schmidhuber, 1997) in the forwards and backwards directions, each of which, like profile HMMs, are strongly linear, left-right models (Eddy, 1998). The choice of BiLSTM was made deliberately to introduce as few changes as possible, ensuring that the observed performance improvements could be attributed primarily to the use of the protein language model, rather than architectural differences from profile HMMs. The output from the BiLSTM layer is subsequently decoded using a stack of three dense layers, scaled to the number of clans or family labels, to produce logits across the prediction space. Probabilities are computed by applying softmax to the logits generated by each model.

## 4 BENCHMARK

In many machine learning contexts, data samples are often assumed to be independent instances drawn from a distribution of the data, justifying random training-test splits. However, this assumption does not hold for sequences in protein domain families, which share evolutionary relationships. Random data splits for protein sequences may lead to performance overestimation, motivating the need to explicitly consider sequence similarity when partitioning data into distinct training and test sets (Söding & Remmert, 2011; Walsh et al., 2016; Jones, 2019; Walsh et al., 2021; Petti & Eddy, 2022). For MDPH-Bench, we aim to create a benchmark that simulates the challenges posed by

Table 1: MDPH-Bench test subset, training, and validation details

Test splits	0-20%	20-40%	40-60%	60-80%	80-100%	Train	Val
Sequences	4,087	37,319	17,446	8,570	5,864	517,936	5,775
Families	543	2,456	1,952	1,731	1,697	14,811	2,097
Clans	180	414	365	341	319	646	388
Coverage	65.31%	59.74%	58.66%	62.02%	56.71%	60.41%	58.79%
Average PID	18.35%	27.64%	42.45%	58.08%	80.01%	NA	45.08%

the remote homology detection problem, where previously unknown or unannotated sequences with little similarity to the training set are especially difficult to detect and annotate. To address this, the MDPH-Bench test set includes a diverse selection of multi-domain proteins, spanning a wide range of sequence similarity to the training set.

#### 4.1 BENCHMARK CREATION

We begin by collecting the 1.2M “seed domains” from Pfam-A Seed 35.0, a set of curated, representative domains for each domain family that are used to build the 20K Pfam profile HMMs (Mistry et al., 2021). We apply BLUE (Petti & Eddy, 2022), a graph-based sequence splitting algorithm, to partition the seed domains into preliminary training and test sets, defining an edge between two domains as their pairwise percent identity (Appendix A.1.1). We use a PID threshold of 25% to split the seed domains, resulting in 560K training domains and 190K test domains. The remaining 450K seed domains were discarded due to sharing > 25% PID with both test and training sets. We then retrieve the full-length sequences corresponding to these representative training and test domains from UniProt release March 2021, a comprehensive database of 230M protein sequences. This results in 517K training sequences. From the test set, we eliminate duplicate sequences also present in the training set. All sequences across both training and test sets are annotated via the `hmmscan` tool from HMMER (Eddy, 2011) with strict inclusion thresholds (E-value < 0.001, bitscore  $\geq$  30) in order to identify domain hits that constitute a “ground truth”, with special care to nested, contiguous domains, which may escape typical processing methods (Appendix A.1.2). For training and test purposes, family and clan labels are only assigned to ground truth domains. We discard sequences in test that do not contain an annotated ground truth domain represented by at least 20 ground truth domains from the same family in the training set, resulting in 73K test sequences. For each test sequence, we compute the maximum PID between any of its domains to any domain in the training set as a proxy for its distance from the training set (Appendix A.1.3). We partition the test set into five subsets based on this maximum PID (Table 1), and a total of 6K validation sequences are sampled uniformly across the test subsets. Such a partition may result in test sequences that, for example, may be placed in the  $80 < \text{PID} \leq 100$  subset due to a single domain closely related to one in the training set, whereas the test sequence may have several other domains that share significantly lower PID with domains in the training set (this is why the average PID is near the lower bound of the max PID range for many of the test subsets in Table 1). The domain coverage, defined as the average percent of residues in a sequence that are labeled by Pfam domains, is similar across all test subsets.

#### 4.2 ADDRESSING POSSIBLE LEAKAGE

We address the potential for unannotated domains to introduce data leakage across the training and test sets by shuffling all subsequences without family and clan labels in the test sequences to disrupt possible domain structures, preserving  $0^{th}$  order residue-composition (Pearson, 2013; Eddy, 2011). The ground truth for protein sequence annotation is fundamentally unknown, relying on inference rather than complete structural and evolutionary knowledge – this is why we must assume natural sequences contain unannotated true domains that new methods may discover. Since PSALM may be sensitive enough to identify unannotated domains, it is trained with these regions shuffled, to mitigate penalties for “false positives” (with respect to the ground truth annotations). Another source of data leakage may arise from the millions of representative sequences from the UniRef50 database release April 2021 (Suzek et al., 2015) used to train ESM-2. We identify that none of the 4,087 sequences in the  $0 < \text{PID} \leq 20$  test subset were present as representative sequences in this version of UniRef50. However, UniRef50 may contain close homologs to the sequences in this test subset.

Model	# Params		Learning Rate	
	Clan	Family	Clan	Family
PSALM <sub>650</sub>	69M	166M	5e-4	5e-5
PSALM <sub>150</sub>	18M	67M	5e-4	5e-5
PSALM <sub>35</sub>	10M	47M	5e-4	5e-5
PSALM <sub>8</sub>	5M	29M	5e-4	5e-5
PSALM <sub>OH</sub>	56M	153M	1e-4	1e-5

Table 2: Number of parameters and learning rates (LR) for PSALM models.

## 5 RESULTS

### 5.1 BASELINES

We establish two baseline methods for comparison. We use HMMER, the current state-of-the-art protein sequence annotation method, to build profile HMMs from MSAs of the ground truth domains in the training set, denoted as HMMER\*, and use these profiles to annotate the test sequences with `hmmscan`. This allows us to evaluate how a state-of-the-art profile HMM method compares to PSALM when using the same training and testing sets. Additionally, we implement a variant of PSALM, denoted as PSALM<sub>OH</sub>, where one-hot embeddings for each amino acid in a protein sequence are utilized instead of embeddings from the pre-trained protein language model ESM-2. This comparison helps discern whether differences in performance between PSALM and HMMER\* are influenced by ESM-2 or solely by the subsequent neural network architecture. Additionally, we assess the scalability of PSALM by evaluating performance across different ESM-2 model sizes (8M, 35M, 150M, and 650M parameters).

More recent deep learning approaches (e.g., ProtENN) are not included as baselines because they address sequence-level or single-domain classification, which is a different problem that is largely irrelevant to domain- or residue-level annotation. We discuss the difference between the two in our related work section (Section 2.2). Biologists prefer domain-level annotation for many reasons, including avoiding the “transitive catastrophe” (lines 27-33), where unrelated sequences cluster through homologous domains (e.g., protein AB shares homology with BC, BC with CD, and all three cluster; but AB shares no homology with CD). Biologists rely on extensive domain-level annotation resources built on a state of the art of profile HMMs.

### 5.2 IMPLEMENTATION DETAILS

Both PSALM and PSALM-onehot are trained using cross entropy loss over the entire sequence for both family and clan annotations. For training all PSALM+ESM-2 models, we use ADAM optimizer (Kingma & Ba, 2014) with initial learning rate 5e-4 for the clan model and 5e-5 for the family model. These values were selected via hyperparameter tuning from across the following learning rates: [1e-3, 5e-4, 1e-4, 5e-5, 1e-5]. A similar hyperparameter search results in a learning rate of 1e-4 for the PSALM<sub>OH</sub> clan model and 1e-5 for the family model. We employ a learning rate scheduler that reduces the learning rate by a factor of  $\sqrt{10}$  if the validation loss fails to decrease over consecutive epochs with an additional early stopping criterion of 5 epochs with no improvement. The effective batch size is 32,768 tokens.

The number of parameters for all PSALM clan and family models are given in Table 2. All models were trained on four NVIDIA A100 80GB GPUs. To accommodate memory limitations on the GPU, all sequences are truncated to a maximum length of 4096 residues. This truncation strategy only applies to approximately 0.25% of sequences across the training and test sets and does not reflect a model limitation – PSALM can be used to annotate sequences of any length provided enough memory. All procedures from the HMMER tool suite use version 3.4 (Eddy, 2011).

### 5.3 METRICS

Protein sequence databases have vastly more negatives than positives, requiring extremely low (essentially zero) and controllable false positive rate (FPR), as false annotations are amplified and

Table 3: PSALM MDPH-Bench residue-level domain annotation results

PID	Model	Clan				Family			
		TPR	FPR	F1	MCC	TPR	FPR	F1	MCC
0-20%	HMMER*	0.694	0.033	0.819	0.642	0.659	0.033	0.810	0.636
	PSALM <sub>650</sub>	<b>0.944</b>	<b>0.022</b>	<b>0.985</b>	<b>0.957</b>	<b>0.750</b>	<b>0.012</b>	<b>0.978</b>	<b>0.947</b>
	PSALM <sub>150</sub>	0.862	0.133	0.912	0.758	0.621	0.050	0.869	0.730
	PSALM <sub>35</sub>	0.729	0.174	0.847	0.620	0.428	0.071	0.721	0.532
	PSALM <sub>8</sub>	0.589	0.214	0.772	0.488	0.211	0.079	0.463	0.293
	PSALM <sub>OH</sub>	0.490	0.100	0.764	0.559	0.089	0.022	0.236	0.203
20-40%	HMMER*	0.907	0.043	0.941	0.862	<b>0.876</b>	0.043	0.939	0.861
	PSALM <sub>650</sub>	<b>0.966</b>	<b>0.020</b>	<b>0.985</b>	<b>0.964</b>	0.845	<b>0.015</b>	<b>0.982</b>	<b>0.959</b>
	PSALM <sub>150</sub>	0.887	0.092	0.925	0.819	0.727	0.036	0.910	0.817
	PSALM <sub>35</sub>	0.799	0.131	0.873	0.709	0.607	0.056	0.825	0.682
	PSALM <sub>8</sub>	0.636	0.192	0.788	0.553	0.353	0.074	0.623	0.452
	PSALM <sub>OH</sub>	0.516	0.107	0.780	0.602	0.102	0.023	0.282	0.251
40-60%	HMMER*	0.951	0.058	0.957	0.898	0.921	0.058	0.956	0.896
	PSALM <sub>650</sub>	<b>0.977</b>	<b>0.020</b>	<b>0.986</b>	<b>0.966</b>	<b>0.924</b>	<b>0.017</b>	<b>0.984</b>	<b>0.964</b>
	PSALM <sub>150</sub>	0.888	0.058	0.927	0.834	0.806	0.026	0.919	0.832
	PSALM <sub>35</sub>	0.826	0.101	0.882	0.736	0.728	0.049	0.866	0.738
	PSALM <sub>8</sub>	0.704	0.158	0.809	0.598	0.532	0.072	0.741	0.567
	PSALM <sub>OH</sub>	0.666	0.104	0.835	0.671	0.159	0.029	0.430	0.363
60-80%	HMMER*	0.974	0.059	0.971	0.924	0.946	0.059	0.970	0.923
	PSALM <sub>650</sub>	<b>0.984</b>	<b>0.018</b>	<b>0.988</b>	<b>0.970</b>	<b>0.957</b>	<b>0.016</b>	<b>0.988</b>	<b>0.968</b>
	PSALM <sub>150</sub>	0.912	0.058	0.940	0.850	0.859	0.028	0.936	0.852
	PSALM <sub>35</sub>	0.845	0.083	0.900	0.761	0.782	0.045	0.888	0.759
	PSALM <sub>8</sub>	0.728	0.154	0.827	0.609	0.605	0.084	0.778	0.583
	PSALM <sub>OH</sub>	0.788	0.094	0.890	0.745	0.216	0.027	0.573	0.478
80-100%	HMMER*	0.977	0.051	0.972	0.935	0.950	0.051	0.971	0.934
	PSALM <sub>650</sub>	<b>0.981</b>	<b>0.015</b>	<b>0.986</b>	<b>0.969</b>	<b>0.967</b>	<b>0.012</b>	<b>0.986</b>	<b>0.968</b>
	PSALM <sub>150</sub>	0.895	0.049	0.929	0.845	0.851	0.024	0.924	0.845
	PSALM <sub>35</sub>	0.812	0.088	0.872	0.732	0.732	0.046	0.853	0.725
	PSALM <sub>8</sub>	0.711	0.114	0.809	0.624	0.601	0.059	0.761	0.600
	PSALM <sub>OH</sub>	0.877	0.066	0.925	0.836	0.282	0.018	0.709	0.630

propagated to additional sequences by later searches. Methods in this field are typically benchmarked for the sensitivity or true positive rate (TPR) they can achieve at a high specificity (low FPR). We also report the F1 score and Matthews Correlation Coefficient (MCC). Here, FPR is defined as the fraction of true negative residues (shuffled, preserving residue composition) incorrectly identified as homologous to a Pfam protein domain family, and TPR is defined as the fraction of residues in ground truth domains correctly identified.

#### 5.4 EVALUATION

We highlight the key observations from the sequence annotation benchmark (Table 3). PSALM<sub>650</sub> demonstrates superior performance in residue-level domain annotation, accurately annotating a substantial portion of true domain regions while consistently calling fewer false positives compared to HMMER\*. Specifically, PSALM<sub>650</sub> reaches higher TPR, F1 and MCC scores at a lower FPR than HMMER\*, with the single exception being family TPR at the 20-40% max PID range test subset. The performance of PSALM<sub>650</sub> is especially noteworthy in the 0-20% max PID range test subset, which constitutes the most difficult to detect domains in MDPH-Bench, as these sequences share very little max sequence similarity with any domain in the training set – PSALM<sub>650</sub> is much more sensitive and specific than HMMER\*. PSALM<sub>OH</sub>, the baseline to ablate the significance of the ESM-2 contextual residue-level embeddings, learns clan-level residue annotations comparably to PSALM<sub>8</sub> and PSALM<sub>35</sub>, especially at the higher max PID range test subsets, but it performs relatively poorly at the family level though it has a consistently low FPR. Additionally, the scaling

Table 4: Family-only PSALM MDPH-Bench results at a fixed FPR of 0.01

PID	Model	Clan				Family			
		TPR	F1	MCC	AUC	TPR	F1	MCC	AUC
0-20%	HMMER*	0.662	0.796	0.629	0.681	0.636	0.790	0.625	0.649
	PSALM <sub>650</sub>	<b>0.922</b>	<b>0.970</b>	<b>0.920</b>	<b>0.934</b>	<b>0.735</b>	<b>0.942</b>	<b>0.868</b>	<b>0.744</b>
	PSALM_F <sub>650</sub>	0.692	0.821	0.659	0.697	0.628	0.806	0.649	0.630
20-40%	HMMER*	0.797	0.885	0.774	0.896	0.774	0.881	0.771	<b>0.867</b>
	PSALM <sub>650</sub>	<b>0.941</b>	<b>0.970</b>	<b>0.931</b>	<b>0.955</b>	<b>0.811</b>	<b>0.934</b>	<b>0.862</b>	0.830
	PSALM_F <sub>650</sub>	0.779	0.877	0.764	0.780	0.746	0.873	0.760	0.746
40-60%	HMMER*	0.502	0.666	0.530	0.882	0.494	0.662	0.527	0.857
	PSALM <sub>650</sub>	<b>0.953</b>	<b>0.974</b>	<b>0.939</b>	<b>0.968</b>	<b>0.889</b>	<b>0.951</b>	<b>0.894</b>	<b>0.909</b>
	PSALM_F <sub>650</sub>	0.827	0.903	0.807	0.831	0.812	0.901	0.806	0.815
60-80%	HMMER*	0.487	0.653	0.503	0.885	0.474	0.646	0.499	0.861
	PSALM <sub>650</sub>	<b>0.967</b>	<b>0.981</b>	<b>0.952</b>	<b>0.977</b>	<b>0.934</b>	<b>0.970</b>	<b>0.927</b>	<b>0.947</b>
	PSALM_F <sub>650</sub>	0.883	0.936	0.854	0.886	0.872	0.935	0.853	0.874
80-100%	HMMER*	0.525	0.687	0.558	0.891	0.515	0.683	0.555	0.866
	PSALM <sub>650</sub>	<b>0.972</b>	<b>0.983</b>	<b>0.962</b>	<b>0.977</b>	<b>0.961</b>	<b>0.981</b>	<b>0.958</b>	<b>0.964</b>
	PSALM_F <sub>650</sub>	0.892	0.940	0.875	0.892	0.887	0.939	0.875	0.887

analysis across the different ESM-2 model sizes (and PSALM sizes, which scale to the selected ESM-2 model) demonstrates increased performance as model size increases up to PSALM<sub>650</sub>.

We conduct an additional ablation experiment in order to study the effect of predicting clan-level annotations as an interpretable intermediate in PSALM (Table 4). We retrain the PSALM<sub>650</sub> family model without providing any predicted clan annotations and denote this model as PSALM\_F<sub>650</sub>. We compare PSALM\_F<sub>650</sub> to HMMER\* and PSALM<sub>650</sub> at a fixed residue-level FPR of 0.01 (the full family-only PSALM table is given in Appendix A.2). We additionally report normalized AUC at this fixed FPR. It is not possible to report AUC without fixing the FPR across all methods, since it is not feasible for the residue-level FPR to reach 1 for HMMER(\*) and PSALM when no alignment or an explicit “none” class is predicted, respectively. Clan predictions are generated by identifying the clan corresponding to the predicted family. PSALM<sub>650</sub> outperforms HMMER\* and PSALM\_F<sub>650</sub> in almost all metrics across all test subsets, demonstrating the advantage of predicting clan annotations prior to family annotations. The clan-level and family-level performance of PSALM\_F<sub>650</sub> are very similar for a given test subset, which is not the case for PSALM<sub>650</sub>, where the clan-level performance is much higher than the family-level performance, especially at the low max PID range test subsets, which represent sequences more distantly related to the training set. This implies that the erroneous family annotations that PSALM\_F<sub>650</sub> predicts are not members of the correct clan, which is typically not the case for PSALM<sub>650</sub>.

## 6 EXAMPLES

In Figure 3, we compare PSALM annotations with those from InterPro (Paysan-Lafosse et al., 2023), a database integrating domain annotations from various sources, focusing on two protein sequences drawn from the test set. This comparison provides insights into PSALM’s performance, highlighting its ability to annotate multiple domains and its high sensitivity in detecting domains.

PSALM is able to identify multiple domains within a single sequence and accurately annotates the three domains in the multi-domain G2/mitotic-specific cyclin-A protein (Fig. 3 A). While most predicted domain boundaries closely match the database annotations, including the extremely short three amino acid region between the N- and C-terminal cyclin domains, some residues that extend past the edges of the first domain (Cyclin-A N-terminal APC/C binding region) exhibit lower assignment probability. This reflects uncertainty in boundary delineation, a common occurrence across annotation methods, which typically differ by a few residues when identifying exact domain boundaries.

PSALM is highly-sensitive and is able to annotate domains that Pfam profile HMMs miss but other databases (and their models) hit as well as domains missed by all methods, as illustrated by the

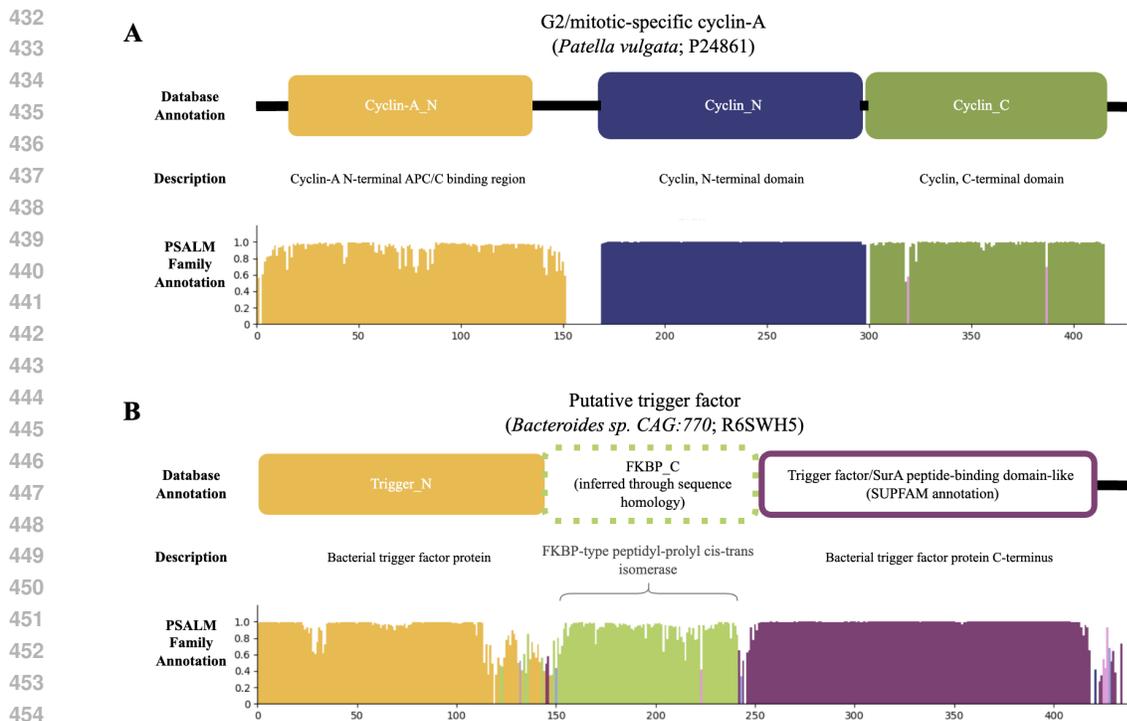


Figure 3: **Comparison of PSALM annotations with those from InterPro for three selected protein sequences.** **A)** PSALM family annotations on a G2/mitotic-specific cyclin-A protein (*Patella vulgata*, P24861) from the 80-100% max PID test subset. **B)** PSALM family annotations on a putative trigger factor protein (*Bacteroides* sp. CAG:770, R6SWH5) from the 20-40% max PID test subset. Filled domains represent InterPro annotations corresponding to Pfam domains. The solid boundary shows a domain from SUPFAM, and the dashed boundary represents a homology-inferred domain not found in InterPro. Approximately 30,000 InterPro proteins share the predicted domain architecture, including trigger factor protein P44837 (*Haemophilus influenzae*).

putative trigger factor protein example (Fig. 3 B). PSALM identifies both N- and C-terminal bacterial trigger factor domains. This is consistent with the integrated InterPro annotation, even though Pfam lacks a C-terminal annotation for this protein.<sup>1</sup> Additionally, PSALM annotates a middle domain that correspond to a FKBP-type isomerase domain, which is missed by all methods for this sequence – we validate this annotation by identifying 30K proteins in InterPro that share this domain architecture (exact order of these three domains) and sequence length.

## 7 CONCLUSIONS

We introduce PSALM, a highly sensitive and specific pLM-based protein sequence annotation method. PSALM extends the capabilities of self-supervised pLMs with just a few hundred thousand protein sequences, enabling interpretable residue-level annotations at both the clan and family levels. Comparisons with InterPro show PSALM’s ability to detect multiple domains, including those currently unannotated. Ablation experiments confirm the importance of pLM embeddings over one-hot encodings and the importance of clan annotations in achieving higher family-level sensitivity and specificity. We find that PSALM performance improves with larger ESM-2 models. We also introduce MDPH-Bench, the first protein sequence benchmark that splits training and test sequences by domain-level sequence similarity for multi-domain proteins. This benchmark minimizes data leakage and enables model evaluation across an evolutionarily-diverse set of proteins.

<sup>1</sup>Putative trigger factor protein R6SWH5 is no longer available in UniProt as of November 7th, 2023. However, the entry and all annotations can be accessed either through UniParc (UniProt Consortium, 2023) or by reannotating with InterProScan (Paysan-Lafosse et al., 2023).

486 Model training and evaluation code as well as sequence IDs for MDPH-Bench are provided in  
487 the supplementary material. Model weights and full sequences are too large to be included in the  
488 supplementary material at the time of submission. After the review period, all code, data, and weights  
489 will be made available online to support reproducibility and protein sequence domain annotation.  
490

## 491 8 LIMITATIONS & FUTURE WORK

492

493 We address the current limitations and potential future research directions of this approach with the  
494 following points.  
495

### 496 8.1 DATA LEAKAGE

497

498 Despite our efforts to mitigate it, information from the test set may still contribute to training through  
499 the millions of sequences used to train ESM-2. While we exclude sequences used to train ESM-2  
500 from our test subset with the lowest maximum PID, homology could still lead to indirect leakage.  
501 However, our maximum PID guarantees help minimize this risk. Ideally, retraining ESM-2 on our  
502 training data would further alleviate this concern. We have not done this in the present work because  
503 of the compute demand for training a pLM like ESM-2, but we hope to rigorously split a much larger  
504 set of proteins to train a pLM from scratch.  
505

### 506 8.2 DOMAIN CALLING

507

508 PSALM cannot distinguish between repeated domains occurring consecutively or accurately resolve  
509 split domains. For example, if a domain repeats immediately after itself, PSALM labels the entire  
510 two domain block instead of recognizing two separate domains within it. Similarly, when a domain  
511 is split, PSALM identifies the two halves as separate domains from the same family, rather than  
512 as originating from a single domain. We aim to address this by explicitly modeling the domain  
513 boundaries and developing domain-calling algorithms.  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

## REFERENCES

- 540  
541  
542 Tristan Bepler and Bonnie Berger. Learning the protein language: Evolution, structure, and function.  
543 *Cell Systems*, 12(6):654–669, 2021.
- 544 Maxwell L Bileschi, David Belanger, Drew H Bryant, Theo Sanderson, Brandon Carter, D Sculley,  
545 Alex Bateman, Mark A DePristo, and Lucy J Colwell. Using deep learning to annotate the protein  
546 universe. *Nature Biotechnology*, 40(6):932–937, 2022.
- 547 Yue Cao and Yang Shen. TALE: Transformer-based protein function Annotation with joint sequence–  
548 Label Embedding. *Bioinformatics*, 37(18):2825–2833, 2021.
- 550 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep  
551 Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 552 Tobias Doerks, Amos Bairoch, and Peer Bork. Protein annotation: detective work for function  
553 prediction. *Trends in Genetics*, 14(6):248–250, 1998.
- 554 Richard Durbin, Sean R Eddy, Anders Krogh, and Graeme Mitchison. *Biological sequence analysis:  
555 Probabilistic models of proteins and nucleic acids*. Cambridge university press, 1998.
- 556 Sean R. Eddy. Profile hidden Markov models. *Bioinformatics (Oxford, England)*, 14(9):755–763,  
557 1998.
- 558 Sean R Eddy. Accelerated Profile HMM Searches. *PLoS Computational Biology*, 7(10):e1002195,  
559 2011.
- 560 Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Yu Wang, Llion Jones, Tom  
561 Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, et al. Prottrans: Toward Understanding  
562 the Language of Life Through Self-Supervised Learning. *IEEE transactions on pattern analysis  
563 and machine intelligence*, 44(10):7112–7127, 2021.
- 564 EMBL-EBI. Protein Classification: What Are Protein Do-  
565 mains? [https://www.ebi.ac.uk/training/online/  
566 courses/protein-classification-intro-ebi-resources/  
567 protein-classification/what-are-protein-domains/](https://www.ebi.ac.uk/training/online/courses/protein-classification-intro-ebi-resources/protein-classification/what-are-protein-domains/), 2024.
- 568 Robert D Finn, Jaina Mistry, Benjamin Schuster-Böckler, Sam Griffiths-Jones, Volker Hollich, Timo  
569 Lassmann, Simon Moxon, Mhairi Marshall, Ajay Khanna, Richard Durbin, et al. Pfam: clans, web  
570 tools and services. *Nucleic Acids Research*, 34(suppl\_1):D247–D251, 2006.
- 571 Tymor Hamamsy, James T Morton, Robert Blackwell, Daniel Berenberg, Nicholas Carriero, Vladimir  
572 Gligorijevic, Charlie EM Strauss, Julia Koehler Leman, Kyunghyun Cho, and Richard Bon-  
573 neau. Protein remote homology detection and structural alignment using deep learning. *Nature  
574 Biotechnology*, pp. 1–11, 2023.
- 575 Michael Heinzinger, Maria Littmann, Ian Sillitoe, Nicola Bordin, Christine Orengo, and Burkhard  
576 Rost. Contrastive learning on protein embeddings enlightens midnight zone. *NAR Genomics and  
577 Bioinformatics*, 4(2):lqac043, 2022.
- 578 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):  
579 1735–1780, 1997.
- 580 Jiajun Hong, Yongchao Luo, Yang Zhang, Junbiao Ying, Weiwei Xue, Tian Xie, Lin Tao, and  
581 Feng Zhu. Protein functional annotation of simultaneously improved stability, accuracy and false  
582 discovery rate achieved by a sequence-based deep learning. *Briefings in Bioinformatics*, 21(4):  
583 1437–1447, 2020.
- 584 David T Jones. Setting the standards for machine learning in biology. *Nature Reviews Molecular  
585 Cell Biology*, 20(11):659–660, 2019.
- 586 Kamil Kaminski, Jan Ludwiczak, Kamil Pawlicki, Vikram Alva, and Stanislaw Dunin-Horkawicz.  
587 pLM-BLAST: distant homology detection based on direct comparison of sequence representations  
588 from protein language models. *Bioinformatics*, 39(10):btad579, 2023.

- 594 Samuel Karlin and Stephen F Altschul. Methods for assessing the statistical significance of molecular  
595 sequence features by using general scoring schemes. *Proceedings of the National Academy of*  
596 *Sciences*, 87(6):2264–2268, 1990.
- 597  
598 Diederik P Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *arXiv preprint*  
599 *arXiv:1412.6980*, 2014.
- 600  
601 Maxat Kulmanov and Robert Hoehndorf. DeepGOPlus: improved protein function prediction from  
602 sequence. *Bioinformatics*, 36(2):422–429, 2020.
- 603  
604 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,  
605 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level  
606 protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- 607  
608 Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, and Alex Rives. Language models  
609 enable zero-shot prediction of the effects of mutations on protein function. *Advances in Neural*  
*Information Processing Systems*, 34:29287–29303, 2021.
- 610  
611 Jaina Mistry, Sara Chuguransky, Lowri Williams, Matloob Qureshi, Gustavo A Salazar, Erik LL  
612 Sonnhammer, Silvio CE Tosatto, Lisanna Paladin, Shriya Raj, Lorna J Richardson, et al. Pfam:  
The protein families database in 2021. *Nucleic Acids Research*, 49(D1):D412–D419, 2021.
- 613  
614 Vamsi Nallapareddy, Nicola Bordin, Ian Sillitoe, Michael Heinzinger, Maria Littmann, Vaishali P  
615 Waman, Neeladri Sen, Burkhard Rost, and Christine Orengo. CATHe: detection of remote homo-  
616 logues for CATH superfamilies using embeddings from protein language models. *Bioinformatics*,  
617 39(1):btad029, 2023.
- 618  
619 Arun Prasad Pandurangan, Jonathan Stahlhacke, Matt E Oates, Ben Smithers, and Julian Gough. The  
620 SUPERFAMILY 2.0 database: a significant proteome update and a new webserver. *Nucleic Acids*  
*Research*, 47(D1):D490–D494, 2019.
- 621  
622 Typhaine Paysan-Lafosse, Matthias Blum, Sara Chuguransky, Tiago Grego, Beatriz Lázaro Pinto,  
623 Gustavo A Salazar, Maxwell L Bileschi, Peer Bork, Alan Bridge, Lucy Colwell, et al. InterPro in  
624 2022. *Nucleic Acids Research*, 51(D1):D418–D427, 2023.
- 625  
626 William R Pearson. Flexible sequence similarity searching with the FASTA3 program package.  
*Bioinformatics Methods and Protocols*, pp. 185–219, 1999.
- 627  
628 William R Pearson. An Introduction to Sequence Similarity (“Homology”) Searching. *Current*  
629 *Protocols in Bioinformatics*, 42(1):3–1, 2013.
- 630  
631 Samantha Petti and Sean R Eddy. Constructing benchmark test sets for biological sequence analysis  
632 using independent set algorithms. *PLOS Computational Biology*, 18(3):e1009492, 2022.
- 633  
634 Roshan Rao, Joshua Meier, Tom Sercu, Sergey Ovchinnikov, and Alexander Rives. Transformer  
635 protein language models are unsupervised structure learners. In *International Conference on*  
*Learning Representations*, 2020.
- 636  
637 Michael Remmert, Andreas Biegert, Andreas Hauser, and Johannes Söding. HHblits: lightning-fast  
638 iterative protein sequence searching by HMM-HMM alignment. *Nature Methods*, 9(2):173–175,  
639 2012.
- 640  
641 Johannes Söding and Michael Remmert. Protein sequence comparison and fold recognition: progress  
and good-practice benchmarking. *Current Opinion in Structural Biology*, 21(3):404–411, 2011.
- 642  
643 Baris E Suzek, Yuqi Wang, Hongzhan Huang, Peter B McGarvey, Cathy H Wu, and UniProt  
644 Consortium. UniRef clusters: a comprehensive and scalable alternative for improving sequence  
645 similarity searches. *Bioinformatics*, 31(6):926–932, 2015.
- 646  
647 Paul D Thomas, Dustin Ebert, Anushya Muruganujan, Tremayne Mushayahama, Laurent-Philippe  
Albou, and Huaiyu Mi. PANTHER: Making genome-scale phylogenetics accessible to all. *Protein*  
*Science*, 31(1):8–22, 2022.

648 UniProt Consortium. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids Research*,  
649 51(D1):D523–D531, 2023.

650

651 Ian Walsh, Gianluca Pollastri, and Silvio CE Tosatto. Correct machine learning on protein sequences:  
652 a peer-reviewing perspective. *Briefings in Bioinformatics*, 17(5):831–840, 2016.

653 Ian Walsh, Dmytro Fishman, Dario Garcia-Gasulla, Tiina Titma, Gianluca Pollastri, Jennifer Harrow,  
654 Fotis E Psomopoulos, and Silvio CE Tosatto. DOME: recommendations for supervised machine  
655 learning validation in biology. *Nature Methods*, 18(10):1122–1127, 2021.

656

657 Ronald J Williams and David Zipser. A Learning Algorithm for Continually Running Fully Recurrent  
658 Neural Networks. *Neural Computation*, 1(2):270–280, 1989.

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

## A APPENDIX

### A.1 BENCHMARK CREATION

#### A.1.1 BLUE

We use the BLUE algorithm (Petti & Eddy, 2022) to split the 1.2M Pfam Seed domains into preliminary train and test sets with a PID threshold of 25%. The pairwise PID between two sequences  $\mathbf{x}$  and  $\mathbf{y}$  is defined as follows:

$$\text{PID}(\mathbf{x}, \mathbf{y}) = \frac{\# \text{ aligned residues}}{\min(\ell(\mathbf{x}), \ell(\mathbf{y}))}, \quad (4)$$

where  $\ell(\mathbf{x})$  and  $\ell(\mathbf{y})$  represent the lengths of sequences  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. If two domains are in the same family, PID is directly calculated from their seed alignment. If two domains are not in the same family but are in the same clan, they are aligned using the `glsearch` tool from the FASTA3 software package (Pearson, 1999), which performs a “global-local” alignment to account for possible large differences in sequence length. If two domains are not in the same clan, they are assumed to share  $< 25\%$  PID.

#### A.1.2 GROUND TRUTH ANNOTATION

We determine “ground truth” by annotating full length sequence with Pfam Seed profile HMMs using `hmmscan` with strict inclusion criteria (E-value  $< 0.001$ , bitscore  $\geq 30$ ). The highest-scoring annotation at each residue is taken as ground truth, but additional post-processing is necessary to ensure that “nested” domain structures are retained. This is accomplished by considering the “match strings” that HMMER generates for an alignment. The match strings contain characters that represent matches, where residues align to a given domain profile, and characters that represent inserts, where unaligned residues are inserted into the sequence relative to the domain profile. Annotating the highest-scoring match state at each residue preserves nested domain structure in the ground truth annotations (Fig. 4).

In a match string, regions with majority matches may contain a few inserts and vice versa. To prevent frequently alternating annotations in the ground truth, we smooth the match and insert states in the

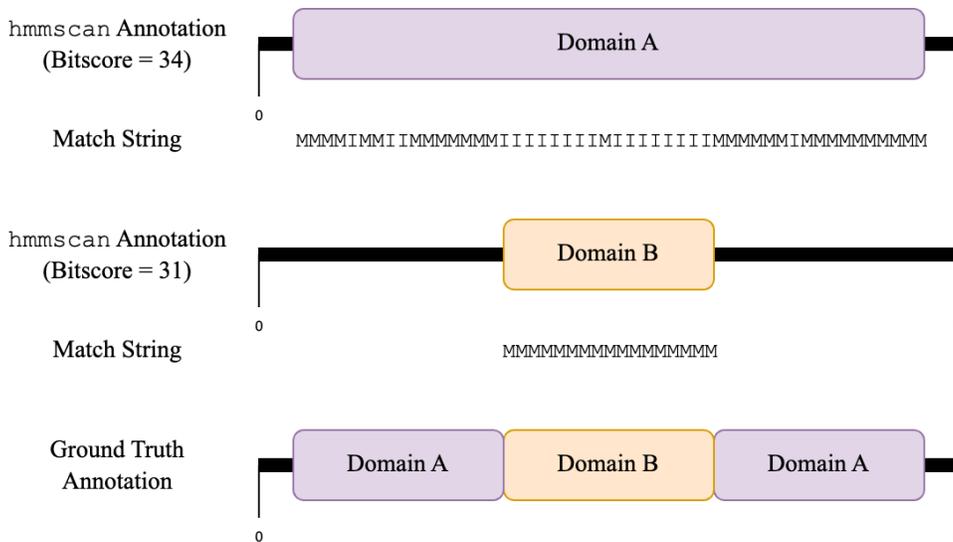


Figure 4: A schematic of nested domains with two domains A and B in the nested format A-B-A. As A is annotated with a higher score than B and overlaps with B, annotating residues only via highest score will fail to include domain B. Using the match state strings to identify smoothed maximal segments preserves the nested domain structure in the ground truth annotation.

match string by identifying maximal scoring segments within the sequence. We assign insert states a positive score and match states a negative score. The segment of the sequence with the greatest aggregate score is known as the maximal segment (Karlin & Altschul, 1990), and all residues in the maximal segment are denoted as insert states. The scores  $s_i$  for each state are inferred from the match string for a given sequence:

$$s_i \propto \log \left( \frac{q_i}{p_i} \right), \quad (5)$$

where  $p_i$  is the frequency with which the state appears in the match string, and  $q_i$  is a state’s target frequency, with  $\sum_i p_i = 1$  and  $\sum_i q_i = 1$ . We set the insert state target frequency at 0.85, the match state target frequency at 0.15, and the length threshold at 20, below which maximal segments are ignored.

### A.1.3 MAXIMUM PID CALCULATION

Once the full length test sequences have been retrieved and subsequently filtered, we compute, for each test sequence, the maximum PID between any of its annotated domains and any annotated domain in train via Algorithm 1. Each test sequence is assigned to a single (out of five) test subset based on its maximum PID.

---

#### Algorithm 1 Percent identity splitting test set

---

**Require:** train sequences  $\mathcal{D}^{tr}$ , test sequences  $\mathcal{D}^{te}$ , Pfam family profile HMMs  $\mathcal{F}$

Initialize an empty dictionary-like structure `record_max_pids`

**for**  $f \in \mathcal{F}$  **do**

`train_domains`  $\leftarrow$  `hmmsearch`  $f$  against  $\mathcal{D}^{tr}$

`test_domains`  $\leftarrow$  `hmmsearch`  $f$  against  $\mathcal{D}^{te}$

**for**  $(\text{domain}, \text{sequence\_id}) \in \text{test\_domains}$  **do**

`MSA`  $\leftarrow$  `hmmalign` `domain` to `train_domains` with  $f$

`domain_pids`  $\leftarrow$  `esl-alipid` `MSA`

`max_pid`  $\leftarrow$  `max`(`domain_pids`)

**if** `sequence_id` not in `record_max_pids` **then**

`record_max_pids`[`sequence_id`]  $\leftarrow$  `max_pid`

**else if** `max_pid` > `record_max_pids`[`sequence_id`] **then**

`record_max_pids`[`sequence_id`]  $\leftarrow$  `max_pid`

**end if**

**end for**

**end for**

Assign each sequence in  $\mathcal{D}^{te}$  to a test split based on max pid

---

The `esl-alipid` tool calculates PID for all pairs of sequences for a given MSA, and is part of the EASEL software package, which can be downloaded together with HMMER (Eddy, 2011).

### A.2 FAMILY-ONLY PSALM FULL RESULTS

We study the effects of predicting clan-level annotations in PSALM by training the PSALM models (across all tested ESM-2 model sizes) without any intermediate predicted clan annotations (Table 5). A subset of this table was included at a fixed FPR of 0.01 in Table 4. Without the intermediate clan predictions, PSALM\_F650 is only competitive with HMMER\* at the clan level for the 0-20% max PID range test subset, though PSALM\_F650 achieves the lowest max FPR across all test subsets.

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

Table 5: Residue-level domain annotation benchmark on Pfam Seed dataset

PID	Model	Clan				Family			
		TPR	FPR	F1	MCC	TPR	FPR	F1	MCC
0-20%	HMMER*	0.694	0.033	<b>0.819</b>	0.642	<b>0.659</b>	0.033	<b>0.810</b>	0.636
	PSALM_F <sub>650</sub>	<b>0.701</b>	<b>0.015</b>	0.827	<b>0.664</b>	0.632	<b>0.015</b>	0.811	<b>0.653</b>
	PSALM_F <sub>150</sub>	0.630	0.034	0.781	0.596	0.540	0.034	0.753	0.576
	PSALM_F <sub>35</sub>	0.560	0.065	0.733	0.523	0.412	0.065	0.670	0.476
	PSALM_F <sub>8</sub>	0.394	0.115	0.599	0.355	0.232	0.115	0.469	0.253
20-40%	HMMER*	<b>0.907</b>	0.043	<b>0.941</b>	<b>0.862</b>	<b>0.876</b>	0.043	<b>0.939</b>	<b>0.861</b>
	PSALM_F <sub>650</sub>	0.781	<b>0.011</b>	0.878	0.764	0.747	<b>0.011</b>	0.873	0.760
	PSALM_F <sub>150</sub>	0.705	0.032	0.833	0.691	0.651	0.032	0.822	0.682
	PSALM_F <sub>35</sub>	0.662	0.058	0.800	0.632	0.581	0.058	0.778	0.614
	PSALM_F <sub>8</sub>	0.479	0.102	0.674	0.462	0.356	0.102	0.605	0.406
40-60%	HMMER*	<b>0.951</b>	0.058	<b>0.957</b>	<b>0.898</b>	<b>0.921</b>	0.058	<b>0.956</b>	<b>0.896</b>
	PSALM_F <sub>650</sub>	0.833	<b>0.012</b>	0.906	0.810	0.816	<b>0.012</b>	0.904	0.809
	PSALM_F <sub>150</sub>	0.785	0.025	0.877	0.759	0.758	0.025	0.873	0.756
	PSALM_F <sub>35</sub>	0.746	0.048	0.846	0.703	0.708	0.048	0.839	0.697
	PSALM_F <sub>8</sub>	0.594	0.087	0.749	0.557	0.520	0.087	0.723	0.535
60-80%	HMMER*	<b>0.974</b>	0.059	<b>0.971</b>	<b>0.924</b>	<b>0.946</b>	0.059	<b>0.970</b>	<b>0.923</b>
	PSALM_F <sub>650</sub>	0.888	<b>0.012</b>	0.938	0.857	0.876	<b>0.012</b>	0.937	0.856
	PSALM_F <sub>150</sub>	0.842	0.025	0.910	0.801	0.823	0.025	0.908	0.799
	PSALM_F <sub>35</sub>	0.792	0.048	0.876	0.733	0.770	0.048	0.872	0.730
	PSALM_F <sub>8</sub>	0.632	0.093	0.776	0.569	0.586	0.093	0.762	0.558
80-100%	HMMER*	<b>0.977</b>	0.051	<b>0.972</b>	<b>0.935</b>	<b>0.950</b>	0.051	<b>0.971</b>	<b>0.934</b>
	PSALM_F <sub>650</sub>	0.892	<b>0.010</b>	0.940	0.875	0.887	<b>0.010</b>	0.939	0.875
	PSALM_F <sub>150</sub>	0.835	0.024	0.903	0.808	0.819	0.024	0.901	0.806
	PSALM_F <sub>35</sub>	0.740	0.047	0.839	0.701	0.720	0.047	0.835	0.697
	PSALM_F <sub>8</sub>	0.661	0.078	0.783	0.609	0.627	0.078	0.774	0.601