



Cite this: DOI: 10.1039/d5dd00472a

Protein language visualizer: a repository for homology exploration with language model embeddings

Javier Espinoza-Herrera,^a María F. Manríquez-García,^b Sofía Medina-Bermejo,^c Ailyn López-Jasso,^d Juan P. Ruiz-Alcocer,^e Adriana Siordia,^a Sarah M. Veskimägi,^a Nate Roethler^a and Adrian Jinich^{*af}

The era of modern AI-driven representations of proteins is here, and moving fast, yet tools for their intuitive visualization and exploration lag behind. Sequence Similarity Networks (SSNs) have long filled this role for alignment-based methods, providing simple but widely adopted platforms for grouping proteins by homology. Building on this foundation, we present the Protein Language Visualizer (PLVis), a modular framework that applies existing pre-trained protein language model (pLM) embeddings, dimensionality reduction, and clustering to generate interactive maps of protein relationships. The central contribution is the PLVis repository, an online resource where thousands of reference proteomes can be compared and annotated through an accessible, interactive interface, much like SSNs became impactful not for their technical novelty but for their broad usability. We first validate that well-separated clusters in PLVis reliably capture homology information, while emphasizing caution when interpreting central “fuzzy” regions. We then illustrate the value of PLVis through case studies spanning individual protein families to full proteome comparisons across *Mycobacterium* and *Plasmodium* species. By combining methodological clarity with broad accessibility, the PLVis repository provides a low-barrier platform for exploring proteomes through the lens of language models.

Received 21st October 2025
Accepted 21st May 2026

DOI: 10.1039/d5dd00472a

rsc.li/digitaldiscovery

Introduction

High-throughput sequencing has accelerated protein discovery, far outpacing functional annotation.¹ UniProtKB now contains over 250 million sequences, yet fewer than 1% are in Swiss-Prot, the manually reviewed section.² Even with automated pipelines, the function of more than 30% of protein-coding genes remains unknown.^{3,4}

Especially in the new era of AI-driven protein representations, where models are powerful but often black-box, visual tools are invaluable for exploring and interpreting large protein collections. Sequence Similarity Networks (SSNs) have long served this role and remain widely adopted.^{5–9} SSNs provide a simple yet effective way to display protein relationships: nodes represent proteins, and edges reflect pairwise similarity scores, which may be scaled or filtered according to the underlying

metric.¹⁰ Tools such as CLANS have long supported the visualization of pairwise sequence similarities, further contributing to the widespread use of network-based representations.¹¹ More recently, web tools like the EFI-EST have made SSNs accessible and popular for probing sequence-function relationships.¹²

Other statistical approaches, such as Profile Hidden Markov Models (HMMs), are also widely used to detect conserved patterns in protein sequences and classify them into families and domains.^{13,14} Originally developed in speech recognition and later adopted in NLP,^{15,16} HMMs underpin many of the curated resources biologists rely on today. However, unlike SSNs, HMM-based results are not typically designed for direct interactive visualization; instead, they are represented through sequence logos, hierarchical trees, or heat maps.^{17–19} Although conceptually distinct, the similarity scores produced by HMM profile searches can also be used in network-based analyses when desired. Complementing these sequence-based strategies, large-scale analyses of structure databases (*e.g.* AlphaFold Protein Structure Database) highlight how structure-based comparisons can also accelerate the annotation of uncharacterized proteins.²⁰

Protein Language Models (pLMs), also adapted from NLP, are the conceptual successors of HMMs.^{21,22} Built on transformer architectures with attention and positional encoding,^{23,24} they are trained through masked-sequence prediction to

^aDepartment of Chemistry and Biochemistry, University of California San Diego, San Diego, USA

^bInstituto Politécnico Nacional, Silao, Mexico

^cUniversidad Autónoma de Baja California, Mexicali, Mexico

^dInstituto Politécnico Nacional, Mexico City, Mexico

^eInstituto Tecnológico y de Estudios Superiores de Occidente, Guadalajara, Mexico

^fSkaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, San Diego, USA. E-mail: ajinich@health.ucsd.edu



learn context-dependent representations of amino acids. The resulting high-dimensional embeddings can represent both individual residues (tokens) and entire protein sequences, providing powerful features for downstream prediction and classification tasks, including structure prediction, as well as protein generation and design.^{23,25–27}

The rise of pLMs and their rich embeddings presents an opportunity to design a new generation of interactive visualization tools for protein similarity. Here, we highlight an accessible approach: interactive exploration of two-dimensional projections of pLM embeddings. Beyond the method itself, our aim is to build a shared resource where such visualizations can be systematically applied across proteomes and protein families. By turning what could be one-off plots into a collective, searchable repository, we hope to make these representations broadly useful for researchers, educators, and even community-driven or citizen science efforts. While several studies have combined protein-language model embeddings with dimensionality reduction or embedding-space visualization to analyze individual protein families or functional subsets (*e.g.*, *via* embedding trees or interactive embedding-space exploration), these efforts remain largely family-centric and do not systematically address species-wide proteome comparison.^{28–31} To our knowledge, no existing resource provides a taxonomically organized, interactive repository for embedding-based visualization across full reference proteomes, integrating clustering, annotation enrichment and cross-species comparative capability.

In this work, we focus on how protein language model embeddings can be explored through 2D projections. Aware of critiques of dimensionality reduction in other fields,

particularly single-cell genomics,^{32–36} we cautiously assess how well these projections preserve sequence relationships. We find that well-separated clusters consistently retain homology information, whereas large, central “fuzzy” regions require caution. Next, through case studies, we illustrate how such projections can support exploratory analysis across scales, from individual protein families to full proteomes, highlighting examples from *Mycobacterium* and *Plasmodium* species. Finally, we introduce the PLVis repository, an interactive resource covering thousands of reference proteomes, alongside a Colab notebook that enables researchers to apply the pipeline to their own data. Together, these contributions establish the PLVis repository as a low-barrier platform for comparative proteomics in the era of protein language models. As with SSNs, the strength of PLVis lies not in methodological novelty but in accessibility and scale: by turning language model embeddings into reusable visual resources, we aim to make them as widely useful for the protein community as SSNs have been for sequence alignments.

Results and discussion

Evaluating PLVis projections: methodological choices, SSN comparison, and functional enrichment

We begin by outlining the PLVis pipeline (Fig. 1). Protein sequences are first embedded with a pLM (*e.g.*, ESM2, ProtT5), then reduced to two dimensions with algorithms such as UMAP, t-SNE, or TriMAP. Clustering methods (*e.g.*, K-means, DBSCAN) identify groups of related proteins, which can be automatically labeled by frequent terms in their annotations. The pipeline is modular, allowing users to mix and match embedding models,

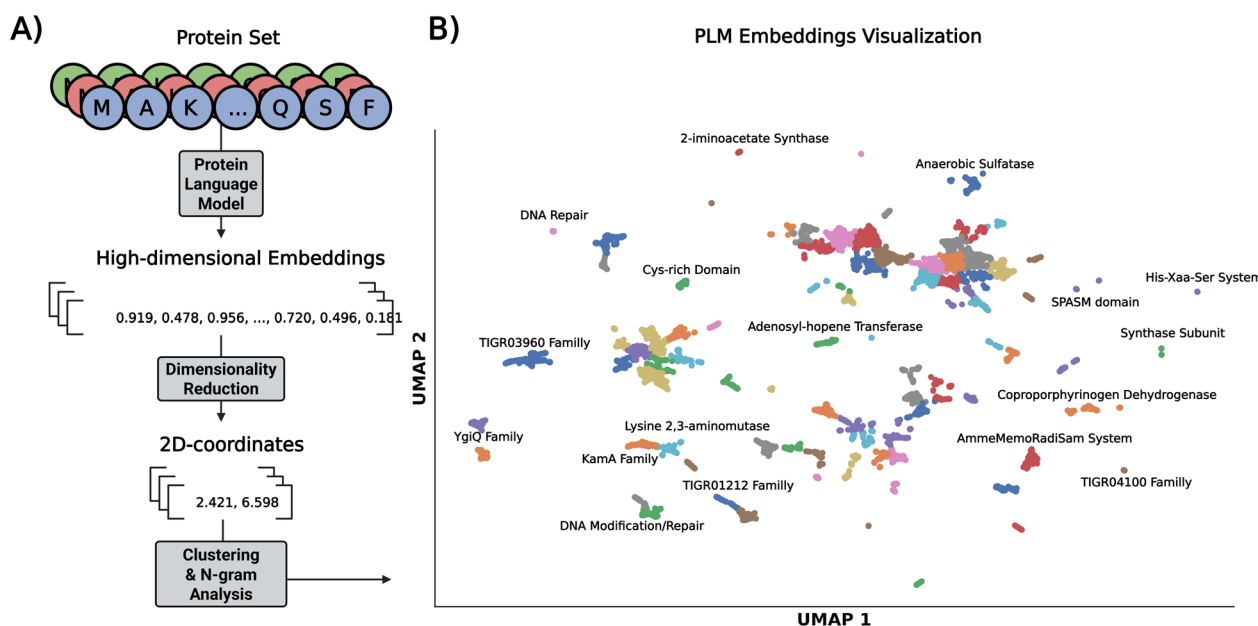


Fig. 1 Schematic overview of the PLVis pipeline. (A) A set of protein sequences is fed to a PLM to obtain embeddings. These embeddings are reduced to two dimensions with a dimensionality reduction algorithm. Next, each point is clustered with its neighbors, and a title that averages the most common words in the protein names is generated using bi-gram analysis. (B) Example visualization of the processed data. The arrows indicate the flow of information through the pipeline, where each step can be performed using the model and reduction algorithm of choice.



dimensionality reduction, and clustering approaches.^{37–41} In the following sections, we use this framework to show how PLVis captures patterns consistent with established homology and orthology annotations, complementing sequence similarity networks and enabling comparative analyses across entire proteomes.

For the analyses in this study, we drew on five datasets: 10 000 radical SAM (rSAM) enzymes, a set of sterol-binding proteins, the full proteome of *M. tuberculosis*, eight *Mycobacterium* proteomes, and five *Plasmodium* proteomes. To compare dimensionality reduction methods, we generated projections with UMAP, t-SNE, and TriMAP under default parameters (Fig. S1 and S2) and assessed clustering quality using the Davies-Bouldin Index (DBI) and Calinski-Harabasz Index (CHI). DBI quantifies cluster quality by comparing within-cluster dispersion to between-cluster separation, whereas CHI measures the ratio of between-cluster to within-cluster dispersion. For the smaller datasets, UMAP consistently produced more compact and well-separated clusters, while for proteome-scale datasets, its performance was intermediate. We further examined UMAP hyperparameters (*e.g.*, *min_dist*, *random seeds*) and found that the overall clustering patterns and case study conclusions remained robust (Fig. S3). Given this balance

of performance and interpretability, we selected UMAP as the default method for subsequent analyses. Lastly, because PLVis is designed as a visualization and exploration tool, clustering is applied to the two-dimensional projections rather than to the original embedding space, ensuring that cluster boundaries correspond to visually separable regions in the final map; this distinction between high-dimensional neighbors and 2D cluster structure is illustrated in Fig. S4. For consistency across datasets, the number of clusters (*K*) was selected by maximizing the silhouette score under default parameters, providing an automated way to identify visually coherent groups.

Next, we compared PLVis to the standard BLAST-based approach, the Sequence Similarity Network (SSN). SSNs are widely used to explore sequence–function relationships by grouping proteins into connected clusters based on pairwise alignment scores.^{9,42–44} A key feature of SSNs is their reliance on a user-defined threshold: at high cutoffs, many proteins appear as isolated nodes, whereas at lower cutoffs more connections are drawn but functional specificity may be lost. To examine how this thresholding compares with PLVis, we analyzed both the 10 000 randomly selected radical SAM (rSAM) enzymes and sterol-binding protein datasets using both approaches (Fig. 2).

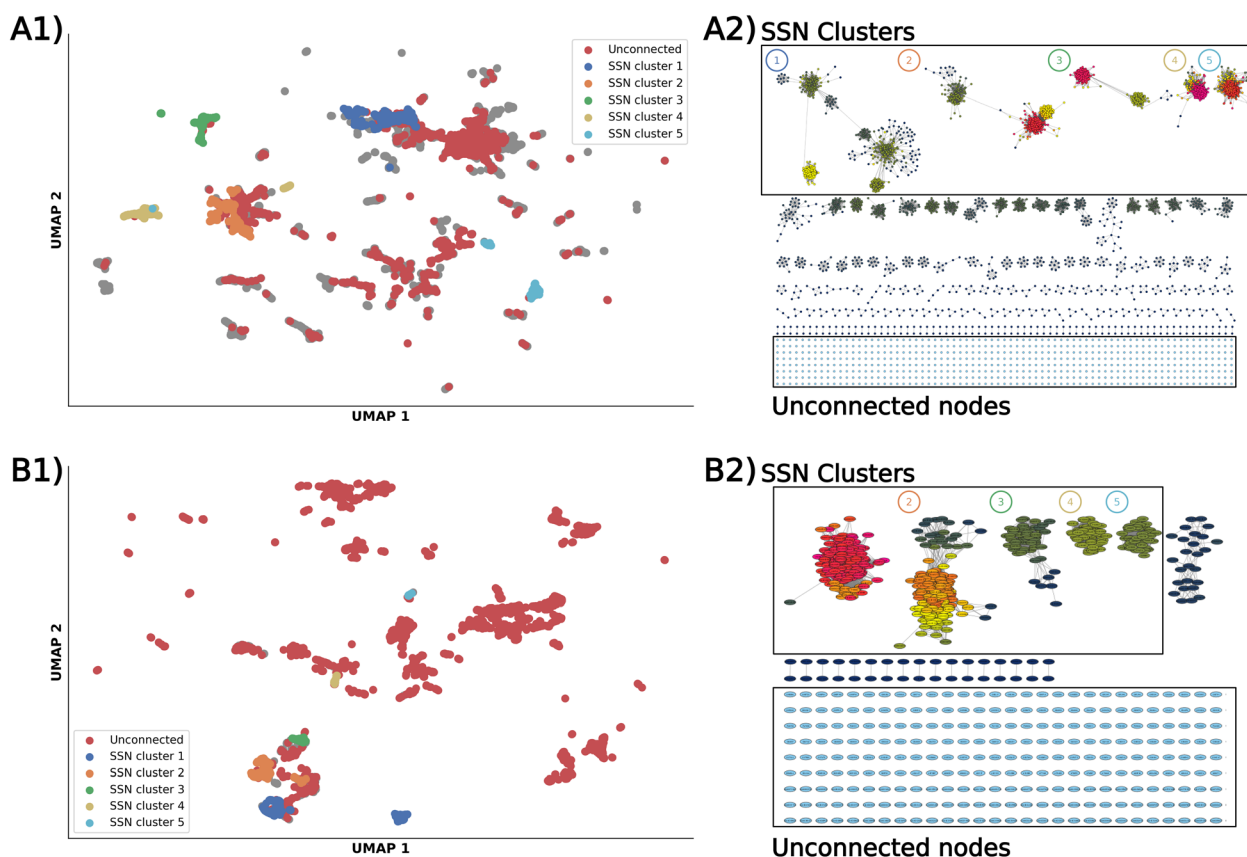


Fig. 2 Comparison of PLVis and corresponding sections of the SSN for the (A) rSAM and (B) sterol datasets. (1) Each color in the PLVis represents a cluster in the SSN (blue: 1, orange: 2, green: 3, yellow: 4, cyan: 5). Proteins colored red in the PLVis represent proteins that are unconnected in the SSN. (2) The SSNs were generated using pairwise sequence alignment scores selected to approximate a 35% sequence identity cutoff. The 5 densest clusters were selected and named accordingly. Nodes situated inside the unconnected region represent proteins that were cut from the threshold.



Fig. 2A1 and B1 shows that the embeddings of both datasets form distinct clusters in two-dimensional space. We selected the five densest clusters from the SSNs (Fig. 2A2 and B2) and examined their correspondence within the PLVis projections, which were partitioned using K-means clustering, yielding 88 clusters for the rSAM dataset and 48 clusters for the sterol dataset. Overall, SSN clusters were largely conserved within the PLVis projections. For instance, SSN cluster 4 in the rSAM comparison was fully contained within PLVis clusters 7 and 57, while SSN clusters 3 and 5 in the sterol-binding comparison were each confined to a single K-means cluster (21 and 32, respectively). A salient characteristic of the SSNs is the large fraction of proteins that appear as isolated nodes at the selected similarity threshold, comprising 1932 proteins (~20% of the rSAM dataset) and 4420 proteins (~80% of the sterol dataset). In contrast, within the PLVis representations, these proteins were redistributed across 75 of the 88 K-means clusters in the rSAM dataset and across all clusters in the sterol dataset, thereby integrating them with related sequences rather than leaving them disconnected.

To evaluate the extent to which PLVis clusters capture biologically meaningful information, we performed enrichment analyses on all K-means clusters using a hypergeometric test for both datasets. For each cluster, the most frequent InterPro annotations (Family, Domain, and Other) were assessed relative to their background frequencies in the full dataset. Among clusters with available annotations, 96% were enriched for an InterPro “Family” term, 74% for “Domain”, and 68% for “Other” in the rSAM dataset, while corresponding values for the sterol dataset were 95%, 90%, and 100%, respectively. Importantly, 93% of the 1932 proteins that appeared as isolated nodes in the rSAM SSN were assigned to PLVis clusters enriched for an InterPro “Family” annotation. For example, PLVis cluster 45 grouped 18 previously unconnected proteins with sequences annotated as belonging to the TatD-associated rSAM family (IPR023821). Likewise, in the sterol-binding representation, PLVis cluster 16 grouped 143 SSN isolated proteins that share a conserved site (IPR020904) within the short-chain dehydrogenase/reductase family (IPR002347). These results illustrate how PLVis can recover meaningful groupings for proteins that SSNs leave disconnected and often excluded from downstream analyses.

To further assess whether PLVis clusters align with curated homology and orthology classifications, we performed systematic enrichment analyses against CATH FunFams and OrthoDB ortholog groups across all datasets (Fig. S5). For each k-means cluster, enrichment was evaluated using a hypergeometric test with Benjamini–Hochberg correction. OrthoDB annotations showed widespread enrichment across clusters in all datasets, indicating strong consistency between embedding-derived groupings and curated ortholog assignments. In contrast, enrichment for CATH FunFams was more variable, reflecting the more limited coverage of structural annotations for large and diverse enzyme families such as rSAM. Together, these results support the use of PLVis as an exploratory framework for organizing proteins in a manner consistent with established homology and orthology resources, without performing explicit

homology or orthology inference; however, the resulting organization is contingent on the underlying embedding model used to represent protein sequences.

PLVis projections preserve local information within well-separated clusters

Aware of the well-documented shortcomings of 2D projections in other biological contexts, we analyzed the properties of PLM embedding projections. In single-cell sequencing studies, dimensionality reduction methods such as t-SNE and UMAP often distort high-dimensional structure. This has been quantified using the Jaccard distance, which measures the overlap between the N nearest neighbors of each point in the original high-dimensional space (ambient space) and in its 2D projection (embedding space). Across 14 single-cell datasets, average Jaccard distances of ~0.7 were reported, indicating substantial loss of local neighborhood information.^{32,36} Such distortions can be considered in terms of local preservation (relationships within clusters), and global preservation (the arrangement of clusters relative to each other in 2D).

We applied this same metric to pLM-based projections to quantify how much clusters preserve neighborhood information. Guided by interactive exploration of annotations across clusters, we hypothesized that well-separated clusters, compact groups that are visually distinct from large, central “fuzzy” regions, more faithfully preserve local relationships in the ambient embedding space, and that this distinction could be quantified rather than assumed. To evaluate this, we analyzed five UMAP datasets, whose properties are summarized in Table S1.

To quantify cluster separation, we calculated silhouette scores for each protein at a fixed number of clusters. We tested thresholds from 0.5 to 0.95, and found that higher cutoffs consistently enriched for clusters with stronger agreement between 2D projections and embedding-space similarity. For clarity, we used 0.95 to illustrate the most stringent case, but the same trend was observed across thresholds (Fig. S6) (shown in blue in Fig. 3). For each protein, we then calculated the Jaccard distance between its neighbors in the ambient embedding space and in the 2D projection. In our analysis (Fig. 3 and S7), well-separated clusters had significantly lower Jaccard distances than other clusters ($p < 0.001$, Mann–Whitney U). As a complementary measure, we compared cosine similarity of high-dimensional embeddings within clusters, which was also significantly higher for well-separated clusters ($p < 0.001$), reflecting that they contain more closely related proteins.

Among the datasets, the eight *Mycobacterium* proteomes showed the clearest effect, with the highest number of well-separated clusters and the strongest agreement between 2D projections and embedding-space similarity. This likely reflects the presence of many conserved protein families shared across closely related species, which are expected to form tight and well-defined groups in embedding space. Importantly, this pattern emerges most clearly in cross-species comparisons, where proteins annotated to the same conserved families cluster together and separate from lineage-specific proteins.



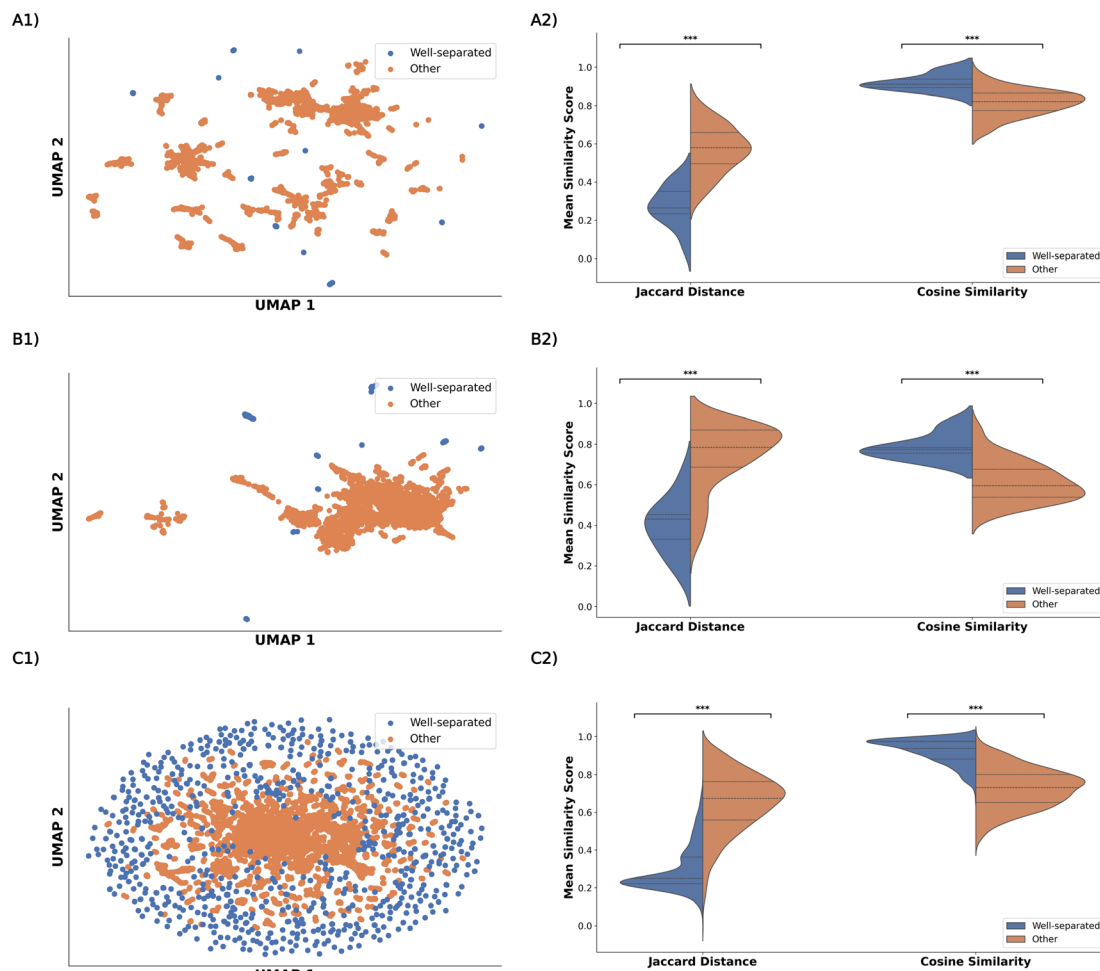


Fig. 3 Well-separated clusters of data are statistically better at conserving high-dimensional data. (1) UMAP plots of PLM embeddings for (A) 10 000 radical SAM enzymes, (B) *M. tuberculosis* proteome, (C) 8 *Mycobacterium* genus proteomes; blue – well-separated clusters, detected by silhouette score above threshold ($S \geq 0.95$), orange – clusters with a silhouette score below the threshold. (2) Violin plots of the average Jaccard distance of proteins and cosine similarity of high-dimensional embeddings within well-separated clusters (blue) and the rest of the clusters (orange). Statistical comparison was performed using the Mann–Whitney U test ($***p < 0.001$).

Such comparative contexts make it particularly straightforward to visually and quantitatively explore conserved protein families across proteomes. Although these trends are expected given the properties of non-linear projections, explicitly quantifying them provides practical guidance for interpreting PLVis visualizations, clarifying when visual separation reflects meaningful embedding-space structure and when it does not.

We then sought to validate the well-established fact that inter-cluster distances in non-linear projections are not particularly meaningful by evaluating whether nearest neighboring clusters have more similar pLM embeddings than randomly selected clusters. Given that non-linear dimensionality reduction techniques like t-SNE and UMAP warp the shape of the data when projecting to lower dimensions, distances between clusters of data points should not be interpreted directly. Using the previously mentioned datasets, we calculated the average cosine similarity for the embeddings of proteins within each cluster and compared it to the inter-cluster cosine similarity with (1) the nearest neighboring cluster and (2) a randomly selected cluster (Fig. 4 and S8).

In Fig. 4 and S8, the violin plots illustrate how cosine similarity varies as we move from proteins within the same cluster to those in the nearest neighboring clusters and finally to random clusters, highlighting trends for both well-separated and poorly-separated clusters. For all datasets, cosine similarity is notably highest within well-separated clusters, aligning with previous observations on local similarity, while poorly-separated clusters show a more gradual decline. We used Cohen's D , an effect size measure that quantifies the magnitude of differences between two distributions (values around 0.2 are considered small, ~ 0.5 medium, and ≥ 0.8 large), to assess two comparisons: (1) intra-cluster *versus* neighboring-cluster similarity scores, and (2) neighboring-cluster *versus* random-cluster similarity scores. These comparisons were performed separately for both well-separated and poorly separated clusters. We also report Mann–Whitney U test results (Table S2), which provide a non-parametric measure of statistical significance; together, the two metrics capture both the size and reliability of the observed effects. When comparing proteins to their nearest neighboring cluster, well-separated clusters showed a sharp decrease in



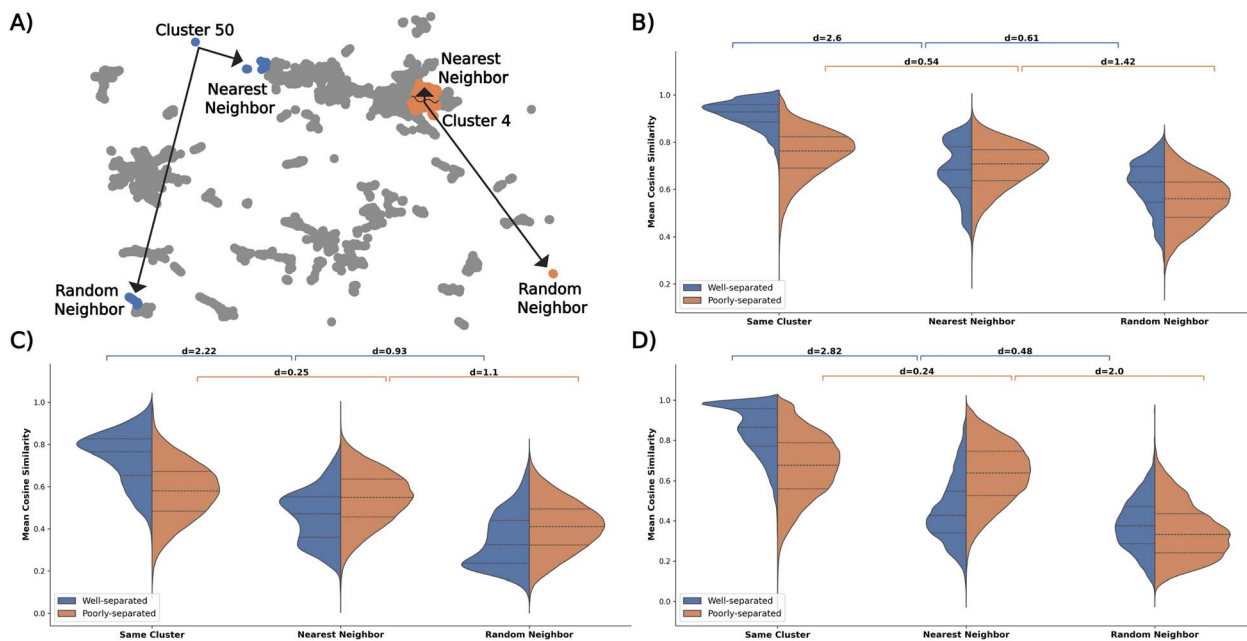


Fig. 4 Distance between clusters in the PLVis projection is not associated with sequence similarity. (A) Example schematic of a well-separated cluster (blue) and a poorly-separated cluster (orange) and their relative positions to their corresponding closest neighboring cluster and a randomly selected cluster; well-separated clusters, detected by silhouette score above threshold ($S \geq 0.95$), poorly-separated clusters with a silhouette score below threshold ($S < 0.5$). (B–D) Violin plots of the mean sequence similarity score for each cluster when comparing its proteins with the nearest neighboring cluster and a randomly selected cluster for (B) 10 000 rSAM enzymes, (C) *M. tuberculosis* proteome, (D) 8 *Mycobacterium* genus proteomes. Significance bars represent the effect size between sets using Cohen's *D*.

similarity relative to within-cluster values, whereas poorly separated clusters showed a more modest decrease (e.g., Cohen's $D = 2.6$ vs. $D = 0.54$, Fig. 4B). This result delineates a clear contrast between regions of the projection where cluster boundaries correspond to embedding-space structure and regions where boundaries are less meaningful. When comparing to a random cluster, the drop in similarity was smaller but more pronounced for poorly separated clusters (e.g., $D = 2.0$ vs. 0.4 , Fig. 4D). This suggests that proteins in poorly separated clusters retain some similarity with nearby clusters in the same cloud. Thus, while absolute distances in 2D should not be overinterpreted, the spatial arrangement of clusters does preserve aspects of the underlying embedding space. While the dimensional reduction serves primarily as a visualization tool, these patterns offer additional context for interpreting both local and global relationships between protein sequences in the visualizations.

In Fig. 4, the violin plots illustrate how cosine similarity varies as we move from proteins within the same cluster to those in the nearest neighboring clusters and finally to random clusters, highlighting trends for both well-separated and poorly-separated clusters. For all three datasets, cosine similarity is notably highest within well-separated clusters, aligning with previous observations on local similarity, while poorly-separated clusters show a more gradual decline. We used Cohen's *D* to measure the effect in two comparisons: (1) between intra-cluster similarity scores and neighboring-cluster similarity scores, and (2) between neighboring-cluster similarity scores and random-cluster similarity scores. These

comparisons were performed separately for both well-separated and poorly-separated clusters. When measuring similarity with the neighboring cluster, proteins belonging to well-separated clusters show a significant drop in the mean, which is not as noticeable when observing the poorly-separated clusters. On the other hand, similar behavior can be observed as we move farther away from the cluster and measure the similarity of proteins with those in a random cluster, but this time, the proteins situated in a poorly-separated cluster show a more significant drop when compared to proteins in well-separated clusters. This implies that sequences in poorly-separated clusters, located in the “fuzzy”, cloud-like aggregation of clusters, share a higher similarity with their surrounding proteins in the cloud-like formation. This pattern suggests that the spatial relationship in the final representation maintains some meaningful reflection of the underlying data structure, even though the absolute distances should not be interpreted directly. While the dimensional reduction serves primarily as a visualization tool, these patterns offer additional context for interpreting both local and global relationships between protein sequences in the visualizations.

PLVis projections reveal conserved protein families across species

Proteins in organisms have evolved to carry out a wide range of biological functions. As species diverge along the phylogenetic tree, their proteomes shift in content and composition. We reasoned that PLVis projections should be particularly useful in



this context, because proteins belonging to conserved families across related species are expected to cluster together in embedding space, making it easier to distinguish broadly conserved groups from lineage-specific proteins. In this section, we focus on two genera of major pathogenic importance, *Mycobacterium* and *Plasmodium*, which include the causative agents of tuberculosis and malaria, respectively.

We first generated a PLM embedding visualization for a subset of species from the genus *Mycobacterium*, a group of over 190 Gram-positive bacterial species belonging to the Actinobacteria phylum. These species range from relatively harmless organisms like *M. smegmatis* to dangerous human pathogens like *M. tuberculosis* and *M. leprae*.^{45,46} These bacteria were traditionally classified by their growth rate (slow or rapid), and recent taxonomic revisions have divided them into five distinct genera: *Mycolicibacterium*, *Mycolicibacter*, *Mycolicibacillus*, *Mycobacteroides*, and *Mycobacterium*.⁴⁷ To demonstrate the value that PLVis projections have in comparing proteomes across organisms, we analyzed and visualized the dataset containing the proteomes of eight *Mycobacterium* species: *M. smegmatis*, *M. fortuitum*, *M. kansasii*, *M. marinum*, *M. leprae*, *M. tuberculosis*, *M. bovis*, and *M. intracellulare* (shown in Fig. 5).

A key insight from visually comparing proteomes across related organisms is the ability to quickly identify which protein families are enriched or expanded in each organism. We thus performed a hypergeometric test with the Benjamini–Hochberg false discovery rate correction to identify the clusters enriched for a single organism. Out of the 1581 k-clusters, 184 (~12%) are enriched and are colored according to their respective organisms in Fig. 5. We found that the three clusters with the lowest

FDR-corrected p -value (clusters 127, 536 & 857) all contained proteins belonging to the PE-Polymorphic GC-Rich (PE-PGRS) family. These proteins are glycine-rich with multiple GGA/GGN repeats and contain a PE domain near the N-terminus of the sequence as well as a high guanine and cytosine (GC) content of approximately 80%.^{48–50} Cluster 857 contains five glycine-rich “uncharacterized” proteins, one of which (A0A7G1IER6) fulfills all previously mentioned qualities (PE domain, GGX motif, and GC content) of a PE-PGRS family protein. Furthermore, all three clusters were not categorized as well-separated, suggesting that they might be closely related to their neighboring clusters, which is further validated by their positions. Both clusters 127 and 536 are close together and linked with cluster 1389, another enriched cluster with PE-PGRS proteins. Cluster 857, although situated on the other side of the projection, is also surrounded by clusters enriched for PE-PGRS family proteins belonging to *M. marinum* (clusters 149 & 1511). This observation is consistent with recent evolutionary analyses showing that the PE-PGRS family is not a homogeneous group but contains subfamilies and specialized members with distinct roles in mycobacterial pathogenesis and host interaction.⁵¹

To understand why the previously mentioned clusters were positioned in separate parts of the visualization, we obtained the AlphaFold structures for proteins belonging to the PE-PGRS clusters and calculated the TM-score between them. The comparison was made between proteins belonging to the same cluster, and random pairs from different clusters were selected as control, shown in the violin plots of Fig. 5. A close look at the plots for each cluster reveals that the visualization separated the protein family according to similarities in their structure, with

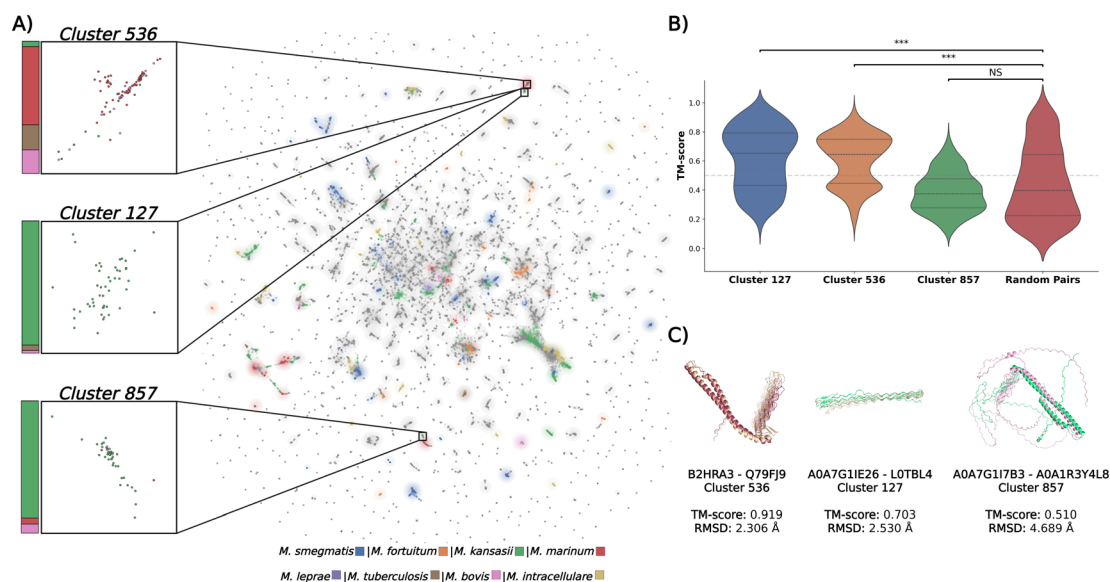


Fig. 5 PLVis for the proteomes of eight *Mycobacterium* species and representative protein structures from three clusters. (A) Clusters in the visualization are colored by enrichment for a particular Mycobacterial species (blue: *M. smegmatis*, orange: *M. fortuitum*, green: *M. kansasii*, red: *M. marinum*, purple: *M. leprae*, brown: *M. tuberculosis*, pink: *M. bovis*, yellow: *M. intracellulare*), clusters in gray are not enriched for a Mycobacterium species. The three most enriched clusters in the projection (127, 536 & 857) are zoomed in, and a color bar showing the fraction of organisms in the cluster is located on their left side. (B) Violin plots of the TM-score (measuring structural alignment) between proteins in clusters 127 (blue), 536 (orange), 857 (green), and random pairs (red). Statistical comparison was performed using the Mann–Whitney U test (** $p < 0.001$) (C) AlphaFold protein structure comparison of proteins in PE-PGRS clusters, colored by organism.



clusters 127 and 536 having most of their scores above the 0.5 threshold. Cluster 857 shows lower scores, but a closer look at the structures in the cluster shows that they have a long disordered region near the C-terminus, which could have impacted the structural comparison. Such structural partitioning aligns with reports that individual PE-PGRS proteins have diverged to acquire specialized functions, suggesting that the separation observed here may reflect true biological heterogeneity within this protein family.⁵¹ We thus infer that the projections can separate proteins belonging to the same family according to their structure, which poses a significant advantage when looking for protein analogs to be used in experimental procedures. However, we reiterate that the distance between both groups of clusters is not a measure of their similarity.

Next, we analyzed the *Plasmodium* genus, consisting of protozoan parasites that require a vertebrate and an invertebrate host to complete their life cycle.⁵² This genus is medically significant as it contains the parasitic species that cause malaria, a vector-borne infection. Five species within this genus are known to infect humans: *P. falciparum*, *P. malariae*, *P. ovale*, *P. vivax*, and *P. knowlesi*.⁵³ Similarly to the previous study, we visualized a dataset containing the proteomes of these five parasites, which is shown in Fig. 6. Compared to the *Mycobacterium* visualization, the *Plasmodium* PLVis has a larger and central poorly-separated/fuzzy region ($S < 0.5$). Of the 1942 k-clusters, approximately 36% were poorly-separated, compared to 14% in the *Mycobacterium* projection.

For this dataset, we repeated the hypergeometric test with the Benjamini-Hochberg false discovery rate correction to

identify clusters enriched for a single organism, which resulted in the identification of 375 (~19%) enriched clusters. We identified 77 enriched clusters that contained proteins exclusively from a single species, a fact that further exemplifies the greater proteomic diversity of this dataset, due to the more complex organisms shown. Because of this greater diversity, one can quickly point out regions in the projection that highlight a specific family of proteins that belong exclusively to a single species in Fig. 6. Such is the case for SICAvAr proteins of *P. knowlesi*, Fam-L proteins of *P. malariae*, and RIFIN proteins belonging to *P. falciparum*.

It has been shown that RIFIN proteins are used by *P. falciparum* to evade the host immune system by binding to immune-inhibitory receptors.⁵⁴ The visualization reveals that most RIFIN proteins are concentrated in three main clusters (38, 1448, and 1522), with only two RIFIN proteins found elsewhere in clusters 1582 and 1666 (A0A143ZXC7 and Q8I209). These outlier proteins are particularly interesting as they are surrounded by members of multiple protein families (REdfSA, tryptophan-rich antigen (TRAGs), and Maurer's clefts two transmembrane (PfMC-2TM) proteins) all of which, including RIFIN proteins, are associated with the infected erythrocyte's membrane.^{55–58} This clustering pattern suggests that A0A143ZXC7 and Q8I209 might also function as erythrocyte surface antigens or membrane proteins. Similarly, the TM-score of these outlier proteins was calculated with the proteins in the RIFIN supercluster and the clusters to which they belong, shown in the violin plots of Fig. 6. We found that the proteins don't share structural similarity with the supercluster, elucidating why they

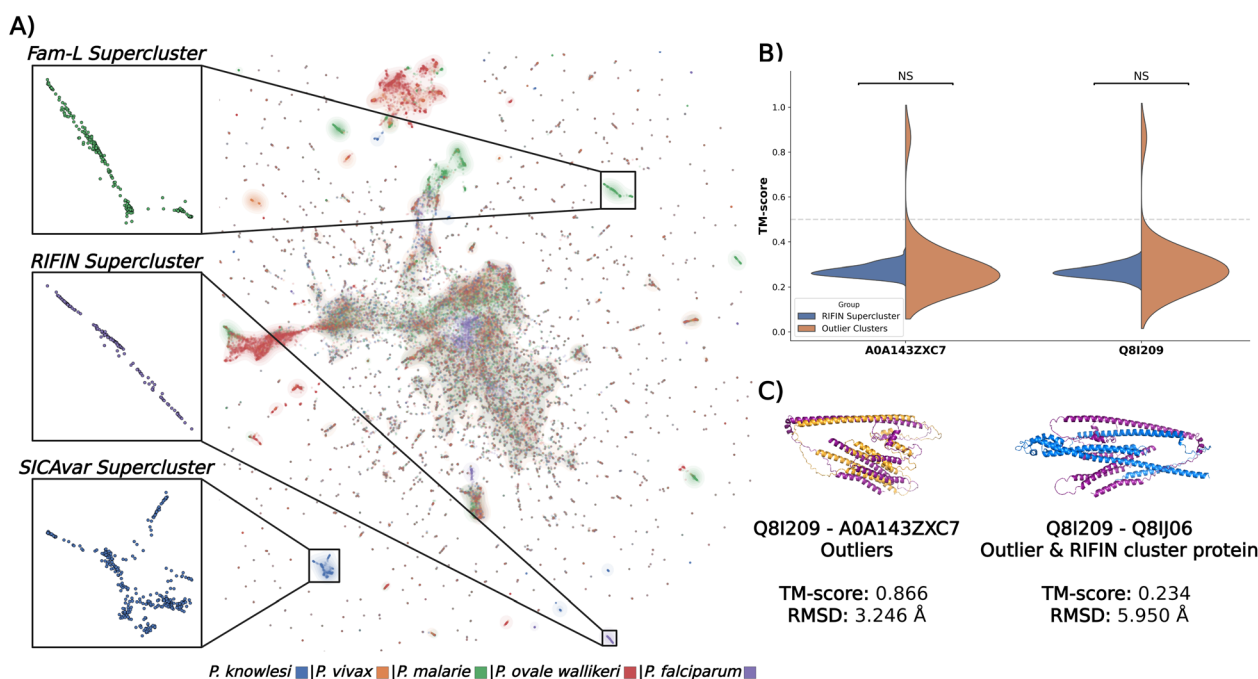


Fig. 6 PLVis for the proteomes of five *Plasmodium* species. (A) Proteins in the visualization are colored according to their species (blue: *P. knowlesi*, orange: *P. vivax*, green: *P. malariae*, red: *P. ovale wallikeri*, purple: *P. falciparum*). The Fam-L, SICAvAr, and RIFIN superclusters are zoomed in. (B) Violin plots of the TM-score between outlier proteins (A0A143ZXC7 and Q8I209) with proteins belonging to the RIFIN supercluster (blue) and outlier clusters (orange). (C) AlphaFold protein structure comparison of Q8I209 (outlier) with A0A143ZXC7 (outlier) and Q8I209 (outlier & RIFIN cluster protein). TM-score: 0.866, RMSD: 3.246 Å; TM-score: 0.234, RMSD: 5.950 Å.



were positioned apart from the other RIFIN proteins. Yet surprisingly, they also don't share significant structural similarity with proteins within each cluster. Nonetheless, their embeddings show relatively high similarity to their corresponding cluster proteins (average cosine similarity ~ 0.7), hinting at a purely functional relation to their neighbors. These observations and associated hypotheses showcase how PLVis can help interactively navigate large-scale protein datasets to reveal biologically significant patterns, while simultaneously providing valuable insights into protein function prediction and pathogen biology.

The PLVis repository, a web portal for comparative proteome analysis within taxonomic families

Having shown the properties of PLVis in the previous sections, we decided to develop an interactive web platform for exploring and comparing UniProtKB reference proteomes, the PLVis Repository. By systematically applying the pipeline across thousands of proteomes and making the results accessible through an interactive interface, the PLVis repository lowers the barrier for engaging with language model representations, whether in research, teaching, or community-driven discovery. The *Mycobacterium* and *Plasmodium* analyses case studies presented above are likewise integrated into the repository as examples of its practical use.

Towards capturing comparative full proteome visualizations across the tree of life, we collected all reference proteomes from UniProtKB. This collection of proteomes covers well-studied model organisms and proteomes of interest in biological research.² Each proteome comparison is performed at the family level, providing a more balanced distribution that offers taxonomic resolution while including enough species for

meaningful comparisons. For each of the available families in UniProtKB, the reference proteomes were retrieved to generate the corresponding PLM embeddings; in the case of outlier taxa with more than ten species, we selected the ten proteomes with the highest BUSCO completeness scores to ensure high-quality and representative comparisons. A total of 4695 reference proteomes are showcased across 3 domains, 3 kingdoms, 67 phyla, 165 classes, 404 orders, 901 families, and 2605 genera.

To facilitate navigation across the different taxonomic groups when first accessing the website, users are greeted with a collapsible tree view that helps them explore the available comparisons. Clicking on each taxon expands the tree, showing the next level of available taxa for each rank, showcasing the relation between the comparisons featured in the repository. The website also features search functionality for specific taxonomic ranks, allowing the rapid retrieval of specific proteome visualizations. If a match is found, the page will either present a list of relevant taxonomic families or redirect users to the associated proteome comparison page if the match corresponds to an available family.

Each page contains a list of the available species at the top, separated by "genus", a PLVis projection of the proteomes, and an enrichment analysis table highlighting overrepresented annotations. Embedding representations for each protein sequence were generated using the ProtT5 language model, followed by dimensionality reduction with the UMAP and t-SNE algorithms. The resulting embeddings were then clustered using K-means to identify structural and functional patterns within each family. Finally, the K-means clusters were named by creating bi-grams of the most frequent words in the protein names present in that cluster.

Each visualization is colored according to the organisms shown for their quick identification, as can be seen in Fig. 7. A

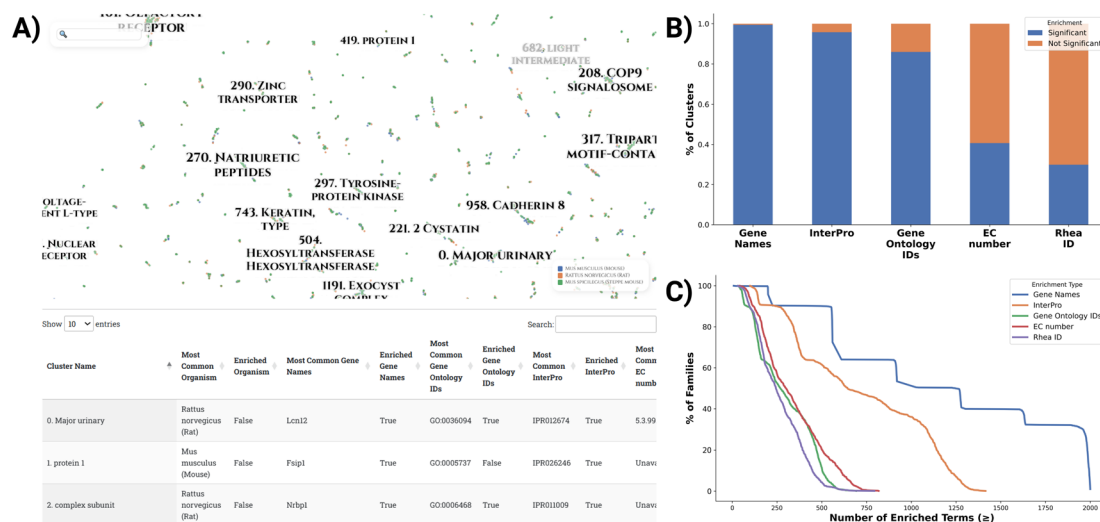


Fig. 7 The PLVis repository shows cluster enrichment of biological significance. (A) Screenshot of the comparison page for the Muridae family. The upper half shows the projection with each protein colored by organism. The bottom half shows a table of enriched terms in each cluster. (B) Stacked bar plots showing the proportion of UMAP clusters enriched for genes, InterPro (IP), gene ontology (GO), EC numbers, and Rhea IDs. Blue bars indicate the fraction of clusters with statistically significant enrichment ($p < 0.05$), while orange bars represent clusters without significant enrichment. (C) Cumulative coverage curves depicting the percentage of UMAP clusters with at least N enriched terms, colored by enrichment type.



search bar located in the upper left corner of each projection allows users to dynamically filter the data by entering keywords. Users can search by UniProt entry ID, gene name, annotation score, and other metadata fields to highlight specific proteins of interest. By examining the visualization and identifying neighboring proteins in the projection space, users can quickly locate proteins with similar embedding representations that, as demonstrated in previous chapters, preserve meaningful functional and structural information.

An enrichment table is located below the visualization to aid users in finding functionally important clusters. Hypergeometric enrichment testing was performed on each cluster by obtaining the most common organism, gene name, InterPro (IP), gene ontology (GO), EC number, and Rhea ID, and comparing their distribution within the cluster against all other clusters. This allows users to detect features that are statistically overrepresented and potentially characteristic of specific protein groups. In Fig. 7 and S9, we show the percentage of clusters that present functional enrichment across all taxonomic families in the PLVis Repository. More than 80% of UMAP and t-SNE clusters in the repository are enriched for gene, IP, and GO information, with more than 50% of the families having more than 500 significant clusters. However, the proportion of enriched clusters is lower for EC number and Rhea ID when compared to the other features. These trends in enrichment emphasize that the protein clusters available are more strongly aligned with functional domain composition than with metabolic reaction identity. Overall, enrichment analysis demonstrates that the clustering captures biologically meaningful groupings, particularly with respect to protein domains (InterPro) and functional annotations (Gene Ontology).

Conclusions

The PLVis pipeline presented here is an efficient and accessible alternative for the visual representation of protein data obtained from PLM embeddings. When used in conjunction with SSNs, these visualizations enhance protein functional annotation by effectively clustering proteins according to their family classifications. For instance, researchers investigating specific protein families and seeking to validate the function of poorly annotated proteins can utilize PLVis projections to rapidly categorize proteins into distinct subfamilies. This clustering facilitates the identification of promising candidates for experimental validation, particularly when minimally annotated proteins (confidence levels 1 or 2) are found in proximity to well-characterized proteins (confidence level 5).

While the primary strength of PLVis lies in its clustering capabilities, it's important to understand both its limitations and flexibility in practical applications. As stated before, due to the limitations of dimensionality reduction, distances in the visualizations aren't meaningful. However, this opens up opportunities for the users to have the liberty to modify cluster coordinates in their datasets, giving meaning to inter-cluster distance based on additional knowledge. For example, clusters can be spatially organized according to various biological

parameters, such as gene expression patterns, protein essentiality profiles, or functional categories (*e.g.*, positioning all redox enzymes in a specific region, or separating transcription factors, transporters, and enzymes). This flexibility in visualization emphasizes the importance of domain expertise and underscores the necessity for users to thoroughly understand both their biological data and the analytical tools at their disposal.

Beyond individual protein analysis and cluster organization, PLVis demonstrates remarkable utility in broader comparative studies. From a biological perspective, PLVis projections demonstrate optimal utility in comparative analyses of complete proteomes across different species. The resultant protein clustering patterns reveal significant biological insights, such as species-specific protein family absences or conserved patterns within taxonomic genera. This approach is particularly valuable for analyzing specific biological relationships, exemplified by host–pathogen interactions, where the visualization can identify clusters of proteins from both organisms that may be implicated in pathogenesis. Such protein clusters provide potential molecular signatures associated with disease mechanisms.

The PLVis repository offers researchers a quick visualization of reference proteome comparisons to find promising protein relationships within taxonomic families. Coupled with the available enrichment tables, the generated projections serve as starting points for deeper biological investigations and hypothesis generation. Expanding the collection of curated case studies through community collaboration could further enhance the website into a better educational and research resource. For this reason, we also make available the link to the PLVis Colab Notebook to assist users in generating their comparisons using the pipeline for the studies above. Together, the PLVis Repository and Colab Notebook provide a scalable platform for visualizing and analyzing large proteomic datasets, helping bridge the gap between massive collections of unannotated proteins and meaningful biological insights.

Author contributions

JEH: conceptualization, data curation, formal analysis, investigation, methodology, software, validation, visualization, writing – original draft. MFMG: data curation, formal analysis. SMB: data curation. ALJ: data curation. JPRA: data curation. AS: data curation. SMV: data curation. NR: data curation. AJ: conceptualization, funding acquisition, resources, supervision, writing – review & editing.

Conflicts of interest

There are no conflicts of interest to declare.

Data availability

Source code for data analysis, along with the proteome data frames needed to replicate our results are available at <https://github.com/javi-eh/PLVis>. The PLVis repository can be accessed at <https://jinichlab.ucsd.edu/plvis/explore.html>. The



PLVis Google Colab referenced in this article is available in the interactive notebook at <https://colab.research.google.com/drive/1s5ug8CYaJ4unJIElxLzcsvxUWPNqWfD?usp=sharing>.

Supplementary information (SI): SI tables and figures. See DOI: <https://doi.org/10.1039/d5dd00472a>.

Acknowledgements

The authors thank the HHMI Hanna Gray Fellowship (GT16787) and the UC San Diego NIH FIRST Grant (NCI 1UA54CA272220) for funding. The authors would also like to thank the ENLACE program at UC San Diego for facilitating the participation of several authors in this research. Figures created in <https://BioRender.com>.

References

- V. de Crécy-lagard, R. Amarin de Hegedus, C. Arighi, J. Babor, A. Bateman, I. Blaby, C. Blaby-Haas, A. J. Bridge, S. K. Burley, S. Cleveland, L. J. Colwell, A. Conesa, C. Dallago, A. Danchin, A. de Waard, A. Deutschbauer, R. Dias, Y. Ding, G. Fang, I. Friedberg, J. Gerlt, J. Goldford, M. Gorelik, B. M. Gyori, C. Henry, G. Hutinet, M. Jaroch, P. D. Karp, L. Kondratova, Z. Lu, A. Marchler-Bauer, M.-J. Martin, C. McWhite, G. D. Moghe, P. Monaghan, A. Morgat, C. J. Mungall, D. A. Natale, W. C. Nelson, S. O'Donoghue, C. Orengo, K. H. O'Toole, P. Radivojac, C. Reed, R. J. Roberts, D. Rodionov, I. A. Rodionova, J. D. Rudolf, L. Saleh, G. Sheynkman, F. Thibaud-Nissen, P. D. Thomas, P. Uetz, D. Vallenet, E. W. Carter, P. R. Weigele, V. Wood, E. M. Wood-Charlson and J. Xu, *Database*, 2022, **2022**, baac062.
- The UniProt Consortium, *Nucleic Acids Res.*, 2025, **53**, D609–D617.
- C. J. Jeffery, *Front. Bioinform.*, 2023, **3**, 1222182.
- K. A. Reynolds, E. Rosa-Molinar, R. E. Ward, H. Zhang, B. R. Urbanowicz and A. M. Settles, *Integr. Comp. Biol.*, 2021, **61**, 2233–2243.
- J. N. Copp, E. Akiva, P. C. Babbitt and N. Tokuriki, *Biochemistry*, 2018, **57**, 4651–4662.
- N. Oberg, T. W. Precord, D. A. Mitchell and J. A. Gerlt, *ACS Bio Med Chem Au*, 2022, **2**, 22–35.
- R. Zallot, N. Oberg and J. A. Gerlt, *Curr. Opin. Biotechnol.*, 2021, **69**, 77.
- A. D. Ashok, J. N. Freitag, I. Irisarri, S. d. Vries and J. d. Vries, *Physiol. Plant.*, 2024, **176**, e14244.
- A. R. Long, E. L. Mortara, B. N. Mendoza, E. C. Fink, F. X. Sacco, M. J. Ciesla and T. M. M. Stack, *Arch. Biochem. Biophys.*, 2024, **757**, 110025.
- H. J. Atkinson, J. H. Morris, T. E. Ferrin and P. C. Babbitt, *PLoS One*, 2009, **4**, e4345.
- T. Frickey and A. Lupas, *Bioinformatics*, 2004, **20**, 3702–3704.
- J. A. Gerlt, J. T. Bouvier, D. B. Davidson, H. J. Imker, B. Sadkhin, D. R. Slater and K. L. Whalen, *Biochim. Biophys. Acta*, 2015, **1854**, 1019.
- S. R. Eddy, *Bioinformatics*, 1998, **14**, 755–763.
- S. C. Potter, A. Luciani, S. R. Eddy, Y. Park, R. Lopez and R. D. Finn, *Nucleic Acids Res.*, 2018, **46**, W200.
- L. Rabiner, *Proc. IEEE*, 1989, **77**, 257–286.
- B. Mor, S. Garhwal and A. Kumar, *Arch. Comput. Methods Eng.*, 2021, **28**, 1429–1448.
- B. Schuster-Böckler, J. Schultz and S. Rahmann, *BMC Bioinf.*, 2004, **5**, 7.
- T. J. Wheeler, J. Clements and R. D. Finn, *BMC Bioinf.*, 2014, **15**, 7.
- A. Krejci, T. R. Hupp, M. Lexa, B. Vojtesek and P. Muller, *Bioinformatics*, 2016, **32**, 9–16.
- I. Barrio-Hernandez, J. Yeo, J. Jänes, M. Mirdita, C. L. M. Gilchrist, T. Wein, M. Varadi, S. Velankar, P. Beltrao and M. Steinegger, *Nature*, 2023, **622**, 637–645.
- D. Ofer, N. Brandes and M. Linial, *Comput. Struct. Biotechnol. J.*, 2021, **19**, 1750.
- T. Bepler and B. Berger, *Cell Syst.*, 2021, **12**, 654.
- A. Madani, B. Krause, E. R. Greene, S. Subramanian, B. P. Mohr, J. M. Holton, J. L. Olmos, C. Xiong, Z. Z. Sun, R. Socher, J. S. Fraser and N. Naik, *Nat. Biotechnol.*, 2023, **41**, 1099–1106.
- A. Chandra, L. Tünnemann, T. Löfstedt and R. Gratz, *eLife*, 2023, **12**, e82819.
- N. Ferruz, S. Schmidt and B. Höcker, *Nat. Commun.*, 2022, **13**, 4348.
- B. Lin, X. Luo, Y. Liu and X. Jin, *Briefings Bioinf.*, 2024, **25**, bbae289.
- R. Schmirler, M. Heinzinger and B. Rost, *Nat. Commun.*, 2024, **15**, 7407.
- M. Heinzinger, M. Littmann, I. Sillitoe, N. Bordin, C. Orengo and B. Rost, *NAR:Genomics Bioinf.*, 2022, **4**, lqac043.
- W. Yeung, Z. Zhou, L. Mathew, N. Gravel, R. Taujale, B. O'Boyle, M. Salcedo, A. Venkat, W. Lanzilotta, S. Li and N. Kannan, *Briefings Bioinf.*, 2023, **24**, bbac619.
- T. Senoner, T. Olenyi, M. Heinzinger, A. Spannagl, G. Bouras, B. Rost and I. Koludarov, *J. Mol. Biol.*, 2025, **437**, 168940.
- A. G. Sangster, C. Dufault, H. Qu, D. Le, J. D. Forman-Kay and A. M. Moses, *PLoS Comput. Biol.*, 2025, **21**, e1012929.
- T. Chari and L. Pachter, *PLoS Comput. Biol.*, 2023, **19**, e1011288.
- D. Kobak and G. C. Linderman, *Nat. Biotechnol.*, 2021, **39**, 156–157.
- A. Diaz-Papkovich, L. Anderson-Trocme, C. Ben-Eghan and S. Gravel, *PLoS Genet.*, 2019, **15**, e1008432.
- A. Dadu, V. K. Satone, R. Kaur, M. J. Koretsky, H. Iwaki, Y. A. Qi, D. M. Ramos, B. Avants, J. Hesterman, R. Gunn, M. R. Cookson, M. E. Ward, A. B. Singleton, R. H. Campbell, M. A. Nalls and F. Faghri, *Patterns*, 2023, **4**, 100741.
- S. Wang, E. D. Sontag and D. A. Lauffenburger, *Cell Syst.*, 2023, **14**, 723.
- A. Elnaggar, M. Heinzinger, C. Dallago, G. Rehawi, Y. Wang, L. Jones, T. Gibbs, T. Feher, C. Angerer, M. Steinegger, D. Bhowmik and B. Rost, *IEEE Trans. Pattern Anal. Mach. Intell.*, 2022, **44**, 7112–7127.
- T. Hayes, R. Rao, H. Akin, N. J. Sofroniew, D. Oktay, Z. Lin, R. Verkuil, V. Q. Tran, J. Deaton, M. Wiggert, R. Badkundri, I. Shafkat, J. Gong, A. Derry, R. S. Molina, N. Thomas, Y. A. Khan, C. Mishra, C. Kim, L. J. Bartie,



- M. Nemeth, P. D. Hsu, T. Sercu, S. Candido and A. Rives, *Science*, 2025, **387**, 850–858.
- 39 L. McInnes, J. Healy, N. Saul and L. Großberger, *J. Open Source Softw.*, 2018, **3**, 861.
- 40 E. Amid and M. K. Warmuth, TriMap: Large-scale Dimensionality Reduction Using Triplets, *arXiv*, 2022, arXiv:1910.00204, DOI: [10.48550/arXiv.1910.00204](https://doi.org/10.48550/arXiv.1910.00204), <https://arxiv.org/abs/1910.00204>.
- 41 Y. Wang, H. Huang, C. Rudin and Y. Shaposhnik, *J. Mach. Learn. Res.*, 2021, **22**, 9129–9201.
- 42 S. Zhao, A. Sakai, X. Zhang, M. W. Vetting, R. Kumar, B. Hillerich, B. San Francisco, J. Solbiati, A. Steves, S. Brown, E. Akiva, A. Barber, R. D. Seidel, P. C. Babbitt, S. C. Almo, J. A. Gerlt and M. P. Jacobson, *eLife*, 2014, **3**, e03275.
- 43 K. H. O'Toole, B. Imperiali and K. N. Allen, *Proc. Natl. Acad. Sci. U. S. A.*, 2021, **118**, e2018289118.
- 44 A. Giorgianni, A. Zenone, L. Sützl, F. Csarman and R. Ludwig, *Microb. Cell Factories*, 2024, **23**, 146.
- 45 E. Tortoli, T. Fedrizzi, C. J. Meehan, A. Trovato, A. Grottola, E. Giacobazzi, G. F. Serpini, S. Tagliazucchi, A. Fabio, C. Bettua, R. Bertorelli, F. Frascaro, V. De Sanctis, M. Pecorari, O. Jousson, N. Segata and D. M. Cirillo, *Infect. Genet. Evol.*, 2017, **56**, 19–25.
- 46 N. L. Bachmann, R. Salamzade, A. L. Manson, R. Whittington, V. Sintchenko, A. M. Earl and B. J. Marais, *Front. Microbiol.*, 2020, **10**, 3019.
- 47 R. S. Gupta, B. Lo and J. Son, *Front. Microbiol.*, 2018, **9**, 67.
- 48 F. De Maio, R. Berisio, R. Manganelli and G. Delogu, *Virulence*, 2020, **11**, 898–915.
- 49 C. D'Souza, U. Kishore and A. G. Tsolaki, *Immunobiology*, 2023, **228**, 152321.
- 50 E. Kramarska, F. De Maio, G. Delogu and R. Berisio, *Biomolecules*, 2023, **13**, 812.
- 51 B. Chen, B. Bajramović, B. Vriesendorp and H. P. Spink, *Biology*, 2025, **14**, 247.
- 52 I. W. Sherman, *Microbiol. Rev.*, 1979, **43**, 453–495.
- 53 S. Antinori, L. Galimberti, L. Milazzo and M. Corbellino, *Mediterr. J. Hematol. Infect. Dis.*, 2012, **4**, e2012013.
- 54 F. Saito, K. Hirayasu, T. Satoh, C. W. Wang, J. Lusingu, T. Arimori, K. Shida, N. M. Q. Palacpac, S. Itagaki, S. Iwanaga, E. Takashima, T. Tsuboi, M. Kohyama, T. Suenaga, M. Colonna, J. Takagi, T. Lavstsen, T. Horii and H. Arase, *Nature*, 2017, **552**, 101–105.
- 55 A. F. Cowman, R. L. Coppel, R. B. Saint, J. Favaloro, P. E. Crewther, H. D. Stahl, A. E. Bianco, G. V. Brown, R. F. Anders and D. J. Kemp, *Mol. Biol. Med.*, 1984, **2**, 207–221.
- 56 B. Wang, F. Lu, Y. Cheng, J.-H. Chen, H.-Y. Jeon, K.-S. Ha, J. Cao, M. H. Nyunt, J.-H. Han, S.-K. Lee, M. P. Kyaw, J. Sattabongkot, E. Takashima, T. Tsuboi and E.-T. Han, *Infect. Immun.*, 2015, **83**, 3083–3095.
- 57 I. Tsarukyanova, J. A. Drazba, H. Fujioka, S. P. Yadav and T. Y. Sam-Yellowe, *Parasitol. Res.*, 2009, **104**, 875–891.
- 58 M. S. Abdel-Latif, K. Dietz, S. Issifou, P. G. Kremsner and M.-Q. Klinkert, *Infect. Immun.*, 2003, **71**, 6229–6233.

