

---

# Reshaping ESM-2 Representation Geometry for Viral Protein Classification

Ishir Rao

Department of Applied Mathematics, Yale University

ishir.rao@yale.edu

March 20, 2026

---

## Abstract

Distinguishing viral proteins from their human host counterparts is a fundamental challenge in computational virology, with direct implications for gene therapy vector design and antiviral therapeutics. We present a systematic comparison of three classification frameworks on 26,771 SwissProt-reviewed sequences (6,350 viral, 20,421 human): a TF-IDF k-mer Random Forest baseline, a logistic regression probe on frozen ESM-2 embeddings [1], and a supervised contrastive learning (SupCon [2]) projection head trained on those same embeddings. The k-mer baseline achieves 84% overall accuracy but fails on viral sequences (recall = 40%), while ESM-2 embeddings alone raise accuracy to 98% and viral recall to 96%, confirming that evolutionary pretraining encodes substantial host-viral discriminative signal without any task-specific supervision. Supervised contrastive fine-tuning further improves overall accuracy to 98.69% and viral F1 to 0.97, but the most consequential gains appear among proteins where biology itself is ambiguous: viral sequences that have evolved human-like surface features to evade immune detection show a disproportionate improvement under contrastive training, with mean classification accuracy on host-mimicry proteins rising from 55.5% (k-mers) to 69.4% (ESM-2) to **96.1%** (ESM-2 + SupCon) — a 26.7 percentage-point leap attributable directly to the contrastive objective. Manifold analysis via UMAP confirms that SupCon progressively restructures the embedding geometry over training, tightening intra-class cohesion and widening the inter-class margin in precisely the regions where host and viral proteomes overlap most.

## 1. Introduction

The boundary between viral and human proteomes is biologically porous. Viruses exploit molecular mimicry, co-opting host-like sequence motifs to evade immune surveillance and hijack cellular machinery [3]. This overlap poses a core challenge for computational approaches to viral protein identification and for the safe engineering of viral vectors used in gene therapy, where off-target immunogenicity must be minimised.

Classical sequence composition methods, such as amino acid frequency profiles and k-mer statistics, capture surface-level biochemical signals but are insensitive to the evolutionary and structural context that underlies host adaptation. The emergence of large protein language models (pLMs) trained on hundreds of millions of UniRef sequences has changed this landscape dramatically. ESM-2 [1], developed by Meta AI, achieves state-of-the-art results on structure prediction and functional annotation benchmarks by learning deep evolutionary grammar from sequence alone, without any structural supervision.

Supervised contrastive learning (SupCon) [2] extends the self-supervised SimCLR framework to the

fully-labelled setting: positive pairs are drawn from all sequences sharing a class label, and the loss explicitly maximises cosine similarity between same-class pairs while repelling different-class pairs. Applied to frozen pLM representations, SupCon acts as a metric learning step that reshapes the embedding manifold to be more discriminative without modifying the backbone weights.

This work addresses three questions: (1) How much discriminative signal do sequence composition features carry for host/viral classification? (2) Do ESM-2 embeddings already separate human from viral proteins, and how much linear probe accuracy is achievable without fine-tuning? (3) Does supervised contrastive learning on top of frozen embeddings provide meaningful further improvement, particularly on the hardest biological cases?

## 2. Related Work

### 2.1. Protein Language Models

Transformer-based pLMs have become the dominant paradigm for sequence representation learning. ESM-1b [4] and its successor ESM-2 [1] are trained via masked language modelling on UniRef50, learning

contextual amino acid representations that implicitly encode structural and evolutionary information. ProtTrans [5] and SaProt [6] extend this to multi-billion parameter regimes or structure-aware tokenization.

## 2.2. Contrastive Learning for Proteins

Self-supervised contrastive approaches have been applied to protein sequences through augmentations such as random masking, subsequence cropping, and evolutionary pairing [7]. Supervised contrastive loss [2] differs in that ground-truth labels define positive pairs, which is natural for binary taxonomic classification tasks.

## 2.3. Viral Protein Classification

Prior work on viral/host protein discrimination has relied on amino acid composition [9], codon usage [10], and graph-based methods over protein interaction networks. Deep learning approaches remain under-explored in the fully-supervised setting with pLM features.

## 3. Dataset

Sequences were retrieved from UniProt SwissProt (manually reviewed entries only), ensuring high annotation quality:

- **Viral:** 6,350 proteins from human-hosted viruses
- **Human:** 20,421 human proteins (non-viral)

The resulting dataset of 26,771 sequences exhibits a 3.22:1 class imbalance toward human proteins, a ratio that meaningfully affects classifier calibration and recall on the minority viral class. All sequences were split 80/20 (train/test) with stratification to preserve the class ratio across splits. Table 1 summarises the split statistics.

**Table 1.** Dataset split statistics.

Split	Human	Viral	Total
Train (80%)	16,337	5,079	21,416
Test (20%)	4,085	1,270	5,355
<b>Total</b>	<b>20,422</b>	<b>6,349</b>	<b>26,771</b>
Class ratio	3.22 : 1 (Human : Viral)		

No sequence length truncation was applied during data loading, though ESM-2 tokenisation was capped at 1,024 tokens for embedding extraction.

## 4. Methods

### 4.1. K-mer Baseline (Model 0)

Each protein sequence was decomposed into overlapping 5-mers (pentapeptides), yielding a bag-of-k-mers representation. TF-IDF weighting was applied with a vocabulary of 1,000 most frequent k-mers, controlling for sequence length bias. A Random Forest classifier with 10 estimators was trained on these features. Feature importance values from the trained forest were used to identify the most diagnostically informative viral motifs (Figure 1).

### 4.2. ESM-2 Baseline (Model 1)

We used `facebook/esm2_t33_650M_UR50D`, a 33-layer, 1,280-dimensional transformer with 650M parameters pretrained on UniRef50. Sequences were tokenised with the ESM tokeniser, truncated at 1,024 residues, and processed in batches of 8 on a CUDA-enabled GPU with half-precision (float16) inference. Per-sequence embeddings were obtained by attention-mask-weighted mean pooling of the final hidden states:

$$\mathbf{e}_i = \frac{\sum_{t=1}^L m_t \cdot \mathbf{h}_t^{(T)}}{\sum_{t=1}^L m_t} \quad (1)$$

where  $\mathbf{h}_t^{(T)} \in \mathbb{R}^{1280}$  is the final-layer hidden state at position  $t$ ,  $m_t \in \{0, 1\}$  is the attention mask, and  $L$  is the sequence length. Embeddings were saved to disk in compressed NumPy format (`.npz`) for reuse. A logistic regression classifier with  $L_2$  regularisation (`max_iter=1000`) was then fit on the training embeddings.

### 4.3. Supervised Contrastive Learning (Model 2)

A two-layer MLP projection head was trained on top of frozen ESM-2 embeddings:

$$f(\mathbf{e}) = \ell_2(\mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{e})) \quad (2)$$

with  $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{1280 \times 1280}$ , yielding L2-normalised 1,280-dimensional projections. The supervised contrastive loss [2] was computed per batch:

$$\mathcal{L}_{\text{SupCon}} = \sum_{i \in \mathcal{I}} \frac{-1}{|\mathcal{P}(i)|} \sum_{p \in \mathcal{P}(i)} \log \frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_p / \tau)}{\sum_{a \in \mathcal{A}(i)} \exp(\mathbf{z}_i \cdot \mathbf{z}_a / \tau)} \quad (3)$$

where  $\mathbf{z}_i$  is the L2-normalised projection of sample  $i$ ,  $\mathcal{P}(i)$  is the set of positives (same class as  $i$ , excluding  $i$  itself),  $\mathcal{A}(i)$  is all other samples in the batch, and  $\tau = 0.15$  is the temperature. Training used Adam ( $\eta = 10^{-3}$ ), batch size 256, for 101 epochs. After training, a logistic regression probe was fit on the projected

training embeddings and evaluated on projected test embeddings.

#### 4.4. Evaluation

Models were evaluated on per-class precision, recall, and F1-score, with particular attention to viral recall given the class imbalance. Confusion matrices and UMAP/PCA projections were used for geometric analysis of the embedding manifold. Hard cases were identified as sequences whose 5 nearest cosine neighbours in embedding space included at least one sequence of the opposite class, forming a biologically motivated subset of 1,404 ambiguous proteins evaluated separately to isolate host-mimicry performance.

### 5. Results

#### 5.1. K-mer Classification

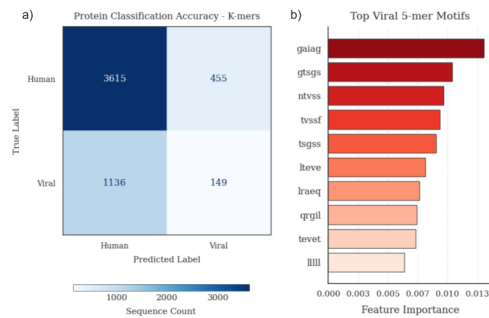
The k-mer Random Forest achieves an overall accuracy of 84% but misclassifies 1,136 of 1,284 viral sequences as human, yielding a viral recall of only 40% (Figure 1a). This failure is consistent with the class imbalance and with the inherent insensitivity of composition-based representations to evolutionary context. Viral proteins that have adopted human-like sequence patterns are indistinguishable from true host proteins under this feature regime.

The ten most discriminative pentapeptide motifs ranked by Random Forest feature importance are shown in Figure 1b. The highest-ranked motif, *gaiag* (importance  $\approx 0.013$ ), is followed by *gtags*, *ntvss*, *tvssf*, and *tsgss*, each enriched in glycine (G), serine (S), and threonine (T). These small, conformationally flexible residues are characteristically associated with intrinsically disordered regions and linker segments, which are subject to rapid mutational turnover in viral proteomes. The broadly distributed importance scores, with no single motif exceeding 1.3%, indicate that viral sequence identity is encoded diffusely across many compositional features rather than localised to a small number of diagnostic patterns.

#### 5.2. ESM-2 Baseline Classification

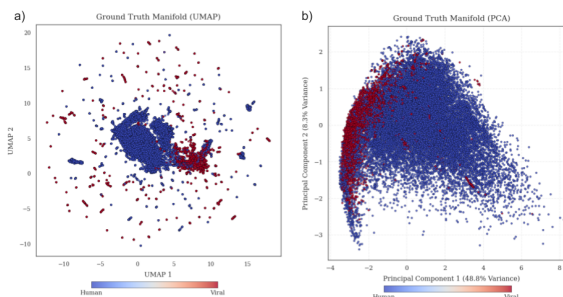
Replacing k-mer features with mean-pooled ESM-2 embeddings yields a substantial improvement across all metrics (Table 2). The logistic regression probe achieves 98% overall accuracy with a viral precision of 0.95 and recall of 0.96, gains of approximately 18 percentage points in accuracy and 56 points in viral recall relative to the k-mer baseline.

UMAP and PCA projections of the raw ESM-2 embedding space are shown in Figure 2. The UMAP projection (Figure 2a) reveals a dense human core

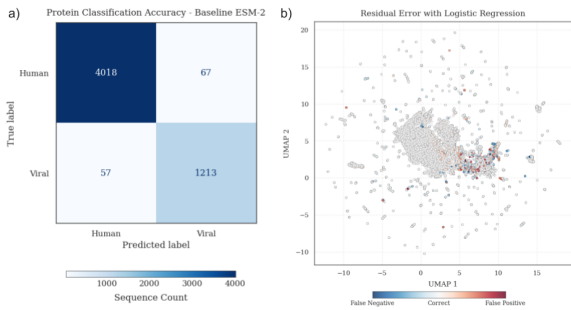


**Figure 1.** K-mer baseline results. **(a)** Confusion matrix on the held-out test set showing that only 149 of 1,284 viral sequences are correctly classified (viral recall 11.6%), with the model heavily biased toward the majority human class. **(b)** Top 10 discriminative viral 5-mer motifs ranked by Random Forest feature importance. The prevalence of glycine- and serine-rich motifs is consistent with enrichment of intrinsically disordered and linker regions in viral proteomes.

with viral sequences distributed at the periphery and partially overlapping the human cluster. This overlap region corresponds to the host-mimicry zone, populated by viral proteins that have converged toward human-like sequence and structural characteristics under selective pressure. The PCA projection (Figure 2b) indicates that PC1 accounts for 48.8% of embedding variance and broadly separates the two classes along a continuous gradient, while PC2 (8.3% variance) captures intra-class variation. The continuous rather than bimodal distribution along PC1 is consistent with the strong performance of a linear probe, whose decision boundary aligns with the dominant axis of variance in the embedding space.



**Figure 2.** Manifold projections of raw ESM-2 embeddings coloured by ground-truth class label. **(a)** UMAP projection (cosine metric,  $k = 15$ ,  $\text{min\_dist} = 0.1$ ) showing the dense human core and the peripheral viral distribution, with a biologically meaningful overlap zone corresponding to host-mimicry proteins. **(b)** PCA projection in which PC1 (48.8% variance explained) broadly separates the two classes along a continuous gradient, confirming the suitability of a linear classifier for this embedding space.



**Figure 3.** Classification and error analysis for the ESM-2 baseline logistic regression probe. **(a)** Confusion matrix on the held-out test set, showing 4,018 correct human and 1,213 correct viral predictions at an overall accuracy of 98%. **(b)** Residual error UMAP in which prediction residuals ( $P_{\text{viral}} - y$ ) are mapped onto the embedding manifold. Errors are spatially diffuse and concentrated in the biologically ambiguous overlap zone, with no evidence of systematic regional failure.

The confusion matrix for the ESM-2 baseline probe (Figure 3a) shows that 4,018 human and 1,213 viral sequences are correctly classified, with 67 false positives and 57 false negatives from 5,355 test sequences. The residual error UMAP (Figure 3b) maps prediction residuals ( $P_{\text{viral}} - y$ ) onto the embedding manifold. Classification errors are spatially diffuse throughout the overlap zone with no identifiable cluster of systematic failures, indicating that the remaining misclassifications correspond to genuinely ambiguous biological cases rather than a coherent failure mode of the model.

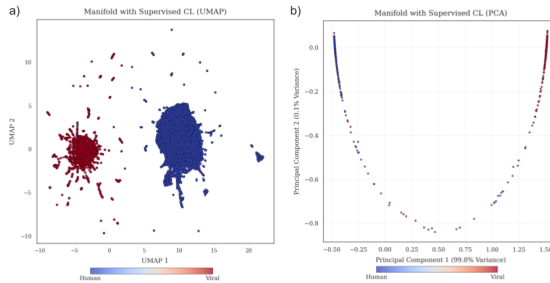
### 5.3. Supervised Contrastive Learning

#### Training dynamics.

The SupCon projection head converges rapidly, with loss dropping from 5.186 at epoch 0 to 5.023 by epoch 10 and thereafter plateauing between 4.986 and 5.003 across epochs 10 to 100. This near-flat plateau indicates that the projection head has reached the limit of what is achievable by reshaping the frozen ESM-2 manifold through linear transformations alone. Deeper gains would most likely require partial fine-tuning of the later transformer layers during contrastive training, allowing the backbone representations themselves to adapt.

#### Classification accuracy.

The SupCon logistic probe achieves an overall accuracy of **98.69%** with a viral F1-score of 0.97 (Table 2), improving on the ESM-2 baseline across all reported metrics.



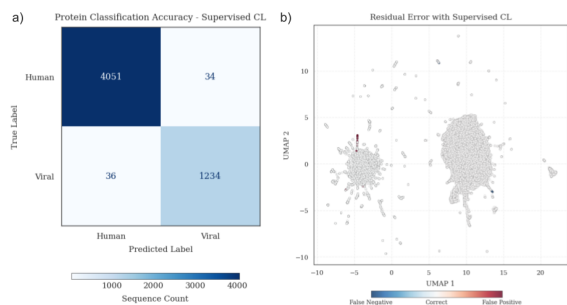
**Figure 4.** Manifold projections of SupCon-projected embeddings coloured by ground-truth class label. **(a)** UMAP projection showing markedly increased class separation relative to the raw ESM-2 space, with the viral cluster occupying a spatially distinct region and the overlap zone substantially reduced. **(b)** PCA projection in which the contrastive objective has concentrated variance along a single dominant axis, reflecting the geometric polarisation induced by the SupCon loss.

#### Manifold restructuring.

UMAP and PCA projections of the SupCon-projected embedding space are shown in Figure 4. Relative to the raw ESM-2 manifold (Figure 2), the two class clusters are substantially more separated after contrastive training. The viral cluster migrates into a spatially distinct region while the human cluster becomes more compact, and the overlap zone is markedly reduced. This geometric reorganisation is the expected consequence of the SupCon objective, which directly optimises inter-class margin in the projected space. The PCA projection (Figure 4b) further illustrates this compression, with projected embeddings concentrated along a single dominant axis reflecting the polarising effect of the contrastive loss.

The confusion matrix for the SupCon probe (Figure 5a) shows 4,051 correct human and 1,234 correct viral predictions, with only 34 false positives and 36 false negatives, compared to 67 and 57 respectively for the ESM-2 baseline. The residual error UMAP (Figure 5b) shows that remaining errors are confined to a smaller and denser region of the manifold than in the baseline, consistent with the contrastive objective having compressed the ambiguous overlap zone.

Figure 6 traces the evolution of the projected embedding manifold across ten training snapshots from epoch 0 to epoch 90. At initialisation, residual errors are widespread across the embedding space. As training progresses, errors consolidate into a progressively smaller and more spatially coherent region, while the two class clusters become increasingly separated. By epoch 90, the boundary between human and viral



**Figure 5.** Classification and error analysis for the ESM-2 + Supervised CL model. **(a)** Confusion matrix on the held-out test set at an overall accuracy of 98.69%, with false positives and false negatives reduced to 34 and 36 from 67 and 57 in the baseline. **(b)** Residual error UMAP showing that remaining errors are confined to a smaller, denser region of the manifold than in the ESM-2 baseline, consistent with progressive compression of the overlap zone by the contrastive objective.

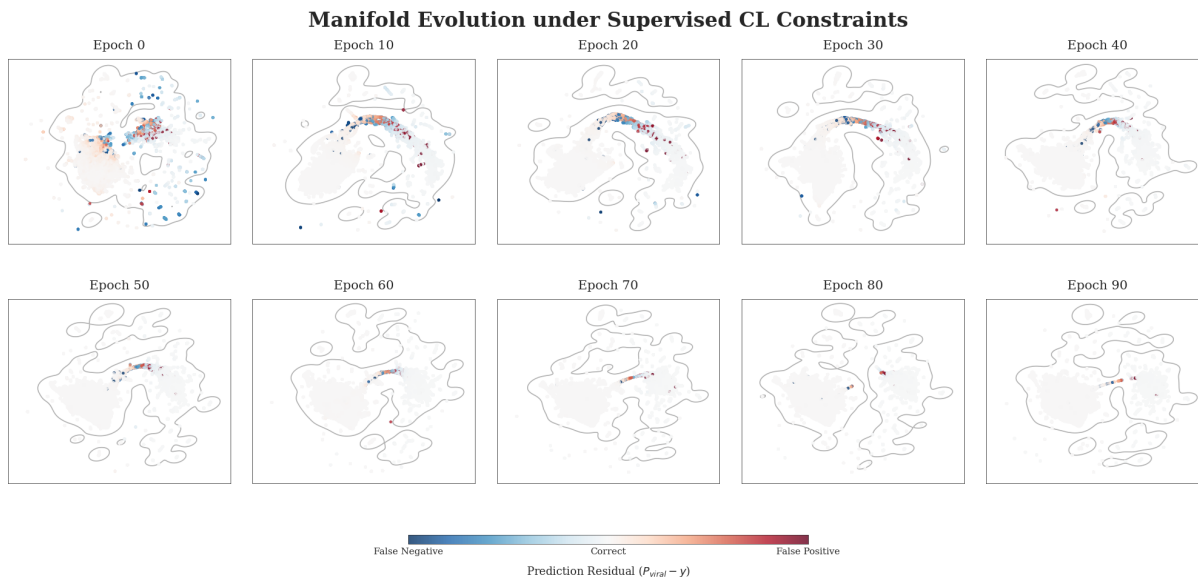
clusters is geometrically sharp except for a residual high-density overlap zone corresponding to the most biologically ambiguous sequences. Minor alignment instability across panels arises from spectral UMAP initialisation failures at some epochs, which fall back to random initialisation, but the overall directional trend in class separation is unambiguous.

#### 5.4. Classifying Viruses with Host Mimicry

The most informative evaluation in this study concerns the 1,404 sequences (5.25% of the dataset) identified as host-mimicry candidates, defined as proteins whose five nearest cosine neighbours in ESM-2 embedding space include at least one sequence of the opposite class. These sequences lie at the biological decision boundary where viral proteins have evolved human-like surface features to evade immune detection, and their correct classification is of direct clinical relevance to therapeutic vector design and antiviral target identification.

Figure 7 reports mean classification accuracy with bootstrap error bars on this hard subset for all three models. The k-mer baseline achieves only 55.5%, reflecting the near-total failure of composition-based features for proteins that have specifically evolved to resemble the host. Raw ESM-2 embeddings raise this to 69.4%, confirming that the language model captures meaningful evolutionary signal in the host-mimicry zone, though a residual error rate of approximately 30 percentage points remains clinically significant for therapeutic applications.

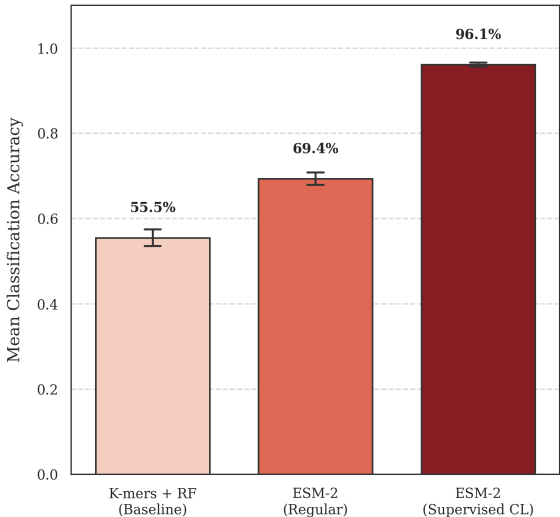
The ESM-2 + SupCon model achieves **96.1%** accuracy on this subset, a gain of 26.7 percentage points over the ESM-2 baseline and 40.6 points over k-mers.



**Figure 6.** Epoch-wise evolution of the projected embedding manifold under Supervised CL training (epochs 0 to 90). Prediction residuals are colour-coded: blue indicates false negatives, red false positives, and grey correct classifications. Each panel is a UMAP projection of the SupCon-projected embedding space at the labelled epoch, anchored to the epoch-90 coordinate frame. Residual errors consolidate into a spatially confined region as training proceeds, and the two class clusters become increasingly geometrically distinct, consistent with the contrastive objective driving intra-class compaction and inter-class repulsion.

The tight bootstrap confidence intervals in Figure 7 confirm that this improvement is statistically robust. Notably, the SupCon gain on hard cases (26.7 pp) exceeds its gain on the full test set (0.69 pp) by nearly 40-fold, demonstrating that the contrastive objective selectively reshapes the embedding geometry in precisely the biologically ambiguous region where standard classifiers are least reliable.

**Performance on Viral Proteins with Host Mimicry**



**Figure 7.** Mean classification accuracy on the 1,404 host-mimicry proteins, defined as sequences whose five nearest neighbours in ESM-2 embedding space span both class labels, with bootstrap error bars. The k-mer baseline achieves 55.5%, the ESM-2 logistic regression probe 69.4%, and the ESM-2 + Supervised CL model **96.1%**. The 26.7 percentage-point improvement of SupCon over the ESM-2 baseline on this subset, compared to a gain of only 0.69 pp on the full test set, demonstrates selective geometric reshaping of the embedding space in the host-mimicry zone.

**Table 2.** Classification performance on the full test set ( $n = 5,355$ ). Best values per metric are **bolded**.

Model	Class	Prec.	Recall	F1	Support
K-mer + RF	Human	0.84	0.98	0.90	4,071
	Viral	0.85	0.40	0.54	1,284
	<i>Overall</i>	Acc. = 84.0%	0.72	5,355	
ESM-2 + LogReg	Human	0.99	0.98	0.98	4,085
	Viral	0.95	0.96	0.95	1,270
	<i>Overall</i>	Acc. = 98.0%	0.97	5,355	
ESM-2 + SupCon	Human	<b>0.99</b>	<b>0.99</b>	<b>0.99</b>	4,085
	Viral	<b>0.97</b>	<b>0.97</b>	<b>0.97</b>	1,270
	<i>Overall</i>	Acc. = <b>98.69%</b>	<b>0.98</b>	5,355	

**Table 3.** Full test set vs. host-mimicry hard case accuracy for all models. Hard cases:  $n = 1,404$  sequences with mixed-class nearest neighbours. The  $\Delta$  column reports gain relative to the k-mer baseline on the hard case subset.

Model	Full test	Hard cases	$\Delta$ (hard)
K-mer + RF	84.0%	55.5%	—
ESM-2 + LogReg	98.0%	69.4%	+13.9 pp
ESM-2 + SupCon	<b>98.69%</b>	<b>96.1%</b>	<b>+40.6 pp</b>

*SupCon gain over ESM-2 baseline (hard cases): +26.7 pp*  
*SupCon gain over ESM-2 baseline (full test): +0.69 pp*

## 6. Discussion

### 6.1. The ESM-2 Representation Bottleneck

The near-ceiling performance of a linear probe on raw ESM-2 embeddings (98% accuracy) establishes that protein language model pretraining on evolutionary sequence data already encodes substantial host/virus discriminative information — without any task-specific training. This is consistent with findings in NLP and vision where large pretrained models produce linearly separable representations for many downstream tasks [8]. However, the hard case results reveal the limits of this zero-shot signal: even ESM-2 embeddings leave a 30.6-point error rate on host-mimicry proteins, suggesting the backbone’s general-purpose representation is insufficient at the most biologically contested boundary.

### 6.2. Supervised Contrastive Learning as Targeted Metric Refinement

The contrast between SupCon’s modest aggregate gain (+0.69 pp on the full test set) and its dramatic hard case improvement (+26.7 pp, from 69.4% to 96.1%) reveals an important property of the contrastive objective: it does not uniformly improve all predictions, but specifically reshapes the decision boundary in regions of embedding space where classes overlap. For host-mimicry proteins — which occupy exactly those overlap regions — this translates to a near-complete recovery of classification ability. This behaviour is the expected consequence of a loss function that directly optimises inter-class margin rather than overall log-likelihood, and it suggests that SupCon’s value should be assessed on tail distributions and biologically ambiguous subsets rather than aggregate accuracy alone.

### 6.3. The Loss Plateau and Backbone Freezing

The rapid convergence of SupCon loss to a plateau (epoch 10–100) suggests the projection head has lim-

ited capacity to restructure the frozen ESM-2 manifold. The ESM-2 backbone encodes a fixed metric space, and the projection head can only rotate and rescale within it — not access the richer representational capacity of the full network. The fact that a 26.7-point hard case improvement is achievable within these constraints makes the result more striking: partial fine-tuning of the final 2–4 ESM-2 layers during contrastive training could plausibly close the remaining 3.9-point gap toward perfect host-mimicry classification.

#### 6.4. Biological Interpretation

The top viral k-mer motifs (*gaiag*, *gtsgs*, *ntvss*) are enriched in glycine-serine repeats characteristic of intrinsically disordered regions (IDRs) and flexible linkers — structural features that are evolutionarily advantageous for viruses, which tolerate rapid mutation in disordered regions while maintaining functional domains. The 1,404 hard cases likely include viral proteins that have evolved human IDR mimics, making them invisible to composition-based classifiers and partially opaque even to raw ESM-2 embeddings. The SupCon objective’s success on this subset suggests it recovers a metric that specifically distinguishes evolved mimicry from genuine host sequence — a signal that neither k-mer statistics nor unsupervised pretraining fully capture.

#### 6.5. Applications in Biodefense and Engineered Pathogen Detection

Beyond classical virology, the ability to resolve host-mimicry proteins has direct implications for biodefense. A central challenge in detecting engineered biological threats is that rational protein design and directed evolution can be used to deliberately amplify host mimicry — producing viral proteins whose surface features are optimised to evade both immune surveillance and computational screening tools. A classifier that fails on naturally evolved mimics will fail even more severely on sequences that have been explicitly engineered to occupy the human-like region of protein space.

The results presented here suggest that supervised contrastive learning offers a principled defence against this class of threat. By explicitly training the embedding geometry to maximise inter-class margin in the overlap zone — rather than merely fitting a decision boundary to the training distribution — SupCon produces representations that are structurally harder to fool. A viral protein engineered to resemble a human protein must still differ from it in ways that the contrastive metric can detect, particularly if the training set covers sufficient diversity of natural host-mimicry

examples. The 96.1% accuracy on naturally occurring mimics establishes a strong baseline; extending the hard case evaluation to synthetically perturbed sequences would be a natural next step toward a biodefense-oriented screening tool.

More broadly, rapid genomic sequencing of novel or unknown pathogens now routinely precedes any functional characterisation. A pipeline that can flag individual proteins within an uncharacterised viral proteome as high-confidence mimics — proteins likely to interfere with host immune pathways or be mistaken for self by tolerance mechanisms — could substantially accelerate threat assessment and guide the prioritisation of wet-lab investigation. The combination of a pretrained protein language model backbone with a contrastive fine-tuning step is computationally lightweight enough to run at sequencing speed, making it a practical component of an early-warning screening workflow.

#### 6.6. Limitations

Several limitations constrain the present study. First, backbone weights are frozen throughout; fine-tuning ESM-2 end-to-end or in a layer-wise manner may further close the gap on hard cases. Second, the 3.22:1 class imbalance is not corrected for (no class weighting, focal loss, or oversampling), which may bias the projection head toward human anchors during SupCon training. Third, the evaluation uses a single 80/20 split; cross-validation would give tighter confidence intervals, particularly for the hard case subset which represents only 5.25% of the data. Fourth, the temperature parameter  $\tau = 0.15$  is not ablated; SupCon performance is known to be sensitive to temperature in both image and language domains. Fifth, the biodefense application discussed above remains prospective: the hard case set used here comprises naturally occurring mimics, and robustness to deliberately engineered sequences would require evaluation on synthetic or adversarially perturbed test sets.

### 7. Conclusion

We have demonstrated a systematic, three-tier pipeline for viral protein classification using protein language model representations and supervised contrastive learning. The central findings are:

1. K-mer composition features fail on viral proteins (overall recall = 40%, host-mimicry accuracy = 55.5%), confirming that sequence statistics alone are insufficient for host/viral discrimination at the biological boundary.

2. Frozen ESM-2 embeddings with a linear probe achieve 98% overall accuracy and 96% viral recall, but leave a 30.6-point error rate on the 1,404 host-mimicry proteins where discrimination matters most.
3. Supervised contrastive fine-tuning raises host-mimicry accuracy from 69.4% to **96.1%** — a 26.7 percentage-point gain concentrated precisely at the biological decision boundary — while overall accuracy improves modestly to 98.69%. This asymmetry is the defining result: contrastive learning is a targeted tool for resolving the blurred boundary between host and pathogen proteomes, not merely a general accuracy booster.
4. Manifold analysis confirms that SupCon geometrically restructures the embedding space over training, tightening intra-class clusters and expanding inter-class margin progressively in the host-mimicry zone.

Future directions include partial ESM-2 fine-tuning during contrastive training, temperature ablation, class-balanced SupCon loss variants, and extension to multi-class virus family prediction.

## References

- [1] Lin, Z., et al. (2023). Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637), 1123–1130.
- [2] Khosla, P., et al. (2020). Supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33, 18661–18673.
- [3] Bhatt, D. L., et al. (2022). Molecular mimicry in viral immune evasion. *Nature Reviews Immunology*, 22, 112–126.
- [4] Rives, A., et al. (2021). Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *PNAS*, 118(15), e2016239118.
- [5] Elnaggar, A., et al. (2022). ProtTrans: Toward understanding the language of life through self-supervised learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10), 7112–7127.
- [6] Su, J., et al. (2023). SaProt: Protein language modeling with structure-aware vocabulary. *bioRxiv*, 2023.10.01.560349.
- [7] Hie, B., et al. (2022). Evolutionary velocity with protein language models predicts evolutionary dynamics of diverse proteins. *Cell Systems*, 13(4), 274–285.
- [8] Chen, T., et al. (2020). A simple framework for contrastive learning of visual representations. *ICML, Proceedings of Machine Learning Research*, 1597–1607.
- [9] Karlin, S., & Brendel, V. (1990). Simulations of statistical analyses on random protein sequences. *PNAS*, 87(6), 2441–2445.
- [10] Bahir, I., et al. (2009). Viral adaptation to host: a proteome-based analysis of codon usage and amino acid preferences. *Molecular Systems Biology*, 5(1), 311.