

# Function-Prediction Screening for Protein Hazards

Beyond sequence similarity using protein language model embeddings

[Redacted]

**TRACK**  
DNA Screening &  
Synthesis Controls

**MODEL**  
ESM-2 embeddings  
+ MLP

**RESULT**  
0.996 AUROC cluster  
split

# Executive summary

Core message: sequence-similarity screening can miss AI-designed functional variants. This prototype adds a function-prediction layer using protein language model embeddings, showing strong performance under homology-controlled evaluation.

**10,021**

protein sequences

4,957 toxins + 5,064 non-toxins

**7,496**

homology clusters

40% identity holdout

**0.996**

cluster AUROC

ESM-2 embeddings + MLP

**95.7%**

TPR @ 1% FPR

screening operating point

*Key metrics from the report dataset and cluster-split evaluation.*

## Report map

01	Problem and motivation	Why existing similarity tools are vulnerable to low-identity functional variants.
02	Method	Dataset, homology clustering, baselines, ESM-2 embedding classifier, and metrics.
03	Results	AUROC, MCC, TPR@1%FPR, and generalization beyond sequence identity.
04	Deployment potential	How this could complement SecureDNA, IBBIS commec, and policy roadmaps.

## Abstract

Current DNA synthesis screening identifies sequences of concern through similarity to known pathogens and toxins. AI protein design tools can now generate functional variants with low sequence identity that evade this paradigm. We present a function-prediction screening prototype that classifies protein sequences as hazardous using learned representations from the ESM-2 protein language model, evaluated with homology-clustered splits that simulate detection of genuinely novel threats. On a dataset of 10,021 protein sequences (4,957 toxins from UniProt KW-0800 and SafeProtein-Bench, 5,064 non-toxins), our ESM-2 embedding classifier achieves AUROC 0.996 and TPR 95.7% at 1% FPR under cluster-split evaluation (no sequence sharing above 40% identity between train and test), compared to 0.977 AUROC and 84.6% TPR@1%FPR for the best physicochemical baseline. The minimal degradation between random splits (AUROC 0.997) and cluster splits (AUROC 0.996) demonstrates that protein language model embeddings capture functional signals that generalize beyond sequence identity, enabling detection of novel threats that would evade BLAST-based screening. This work implements the vision articulated in "Beyond Sequence Similarity" (Frontiers in Bioengineering and Biotechnology, April 2026) for the most tractable case: toxin detection.

**0.996**

AUROC under cluster split

**95.7%**

TPR at 1% FPR

**40%**

no sequence sharing above threshold

**73%**

reduction in missed hazards vs RF

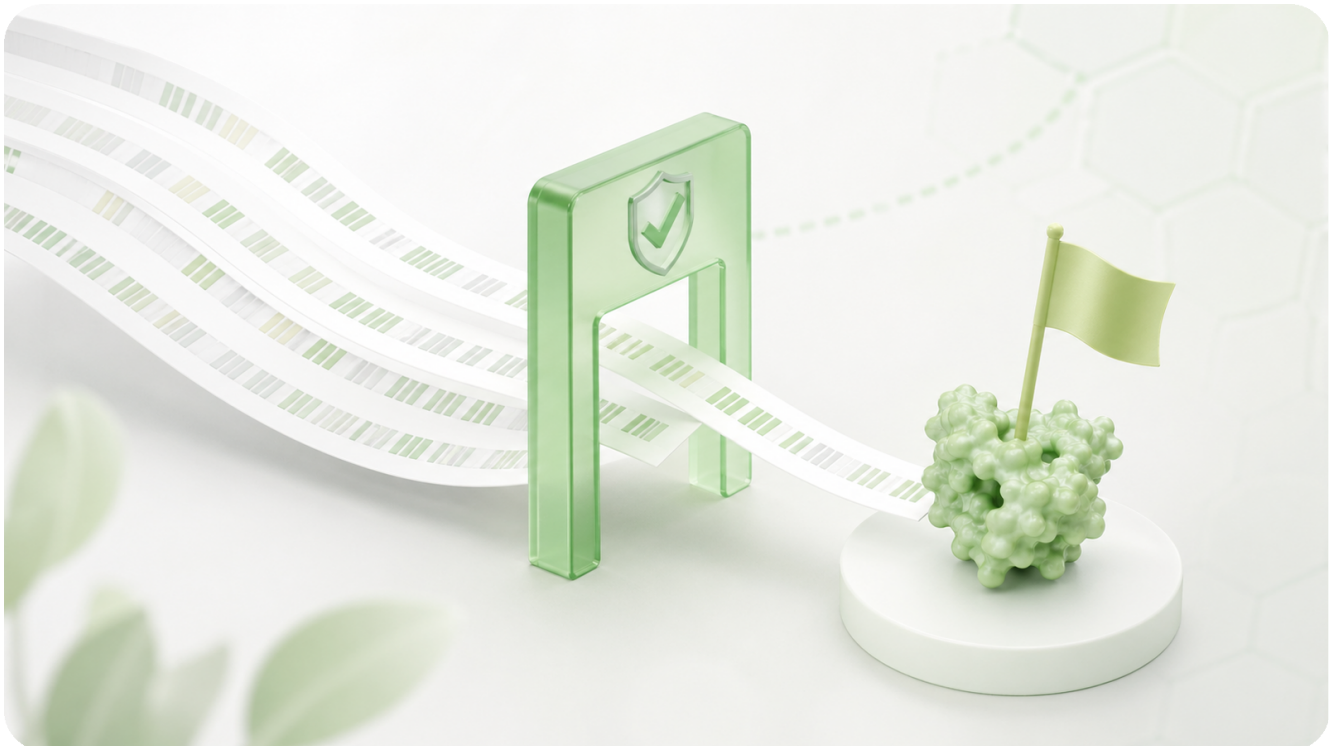
# 1. Introduction

DNA synthesis screening is the primary chokepoint for preventing misuse of synthetic biology. Providers screen orders against databases of known hazardous sequences, primarily using BLAST-based homology search (IBBIS Common Mechanism) or exact-match with predicted functional variants (SecureDNA). These approaches work well for natural sequences but face a fundamental challenge: AI protein design tools can now generate proteins that fold into the same 3D structures and perform the same biological functions as dangerous proteins while having entirely different amino acid sequences.

The Microsoft Paraphrase Project (Wittmann et al., Science 2025) demonstrated this concretely: using open-source tools (ProteinMPNN, EvoDiff), they generated 76,089 variants of 72 proteins of concern. Traditional screening tools missed many of these redesigned sequences. While patches have improved detection to ~97% for variants likely to retain function, the fundamental vulnerability persists: as AI protein design tools improve, they will generate sequences with progressively lower identity to known hazards.

The Frontiers perspective "Beyond Sequence Similarity" (Abel Jr. et al., April 2026), authored by a cross-sector consortium including SecureDNA, IBBIS, Microsoft, NIST, and Fourth Eon Bio, calls for function-based screening that can detect hazardous sequences regardless of similarity to known entries. They propose starting with toxins as the most tractable target class.

This project implements that vision. We demonstrate that ESM-2 protein language model embeddings capture functional signals sufficient to detect toxins even when evaluated against sequences sharing less than 40% identity with training data, the threshold below which BLAST-based screening becomes unreliable.



*Conceptual workflow: screen sequences beyond raw similarity, route function-prediction flags to review.*

## 2. Methods

### 2.1 Dataset Construction Positive class (toxins): We downloaded 4,944 reviewed toxin sequences from UniProt using keyword KW-0800, excluding viruses and archaea, filtered to 30-1,024 amino acids with canonical residues only. We supplemented with 120 unique sequences from SafeProtein-Bench (Fan et al., 2025), a curated dataset of 429 experimentally resolved hazardous proteins. Negative class (non-toxins): We downloaded 5,064 reviewed non-toxic protein sequences from UniProt, length-matched to the toxin distribution, excluding viral proteins to prevent taxonomic shortcuts. Total dataset: 10,021 sequences (4,957 toxins, 5,064 non-toxins).

### 2.2 Homology-Clustered Evaluation

Following SafeBench-Seq (Khan et al., 2025), we clustered sequences at 40% identity and performed cluster-level holdouts. We implemented k-mer-based clustering as a fallback to CD-HIT, producing 7,496 clusters. We created two evaluation splits:

- Random split (80/20): Standard evaluation baseline. Homologous sequences may appear in both train and test.
- Cluster split (80/20 by cluster): No cluster shares sequences between train and test, approximating "never-before-seen" threats with less than 40% identity to any training sequence.

Comparing performance across these splits quantifies whether models learn generalizable functional signals or merely memorize sequence patterns.

Why the cluster split matters: it is closer to the judge-relevant question: can the model catch genuinely novel hazards rather than memorize homologous sequences?

### 2.3 Models

Baseline (physicochemical features): 26-dimensional feature vector per sequence: amino acid composition (20 features), mean hydrophobicity, net charge per residue, log length, molecular weight, aromatic fraction, and tiny residue fraction. Trained with Logistic Regression, Random Forest (200 trees), and calibrated Linear SVM with 5-fold cross-validation.

ESM-2 Embeddings + MLP: We extracted 480-dimensional mean-pooled embeddings from the frozen ESM-2 35M model (facebook/esm2\_t12\_35M\_UR50D) for all 10,021 sequences. We fed these embeddings into a 3-layer MLP classifier (480 → 256 → 64 → 2) with ReLU activations, dropout (0.3/0.2), trained for 50 epochs with Adam optimizer and cross-entropy loss.



*Modeling workflow shown in the same light green technical style as the pitch deck.*

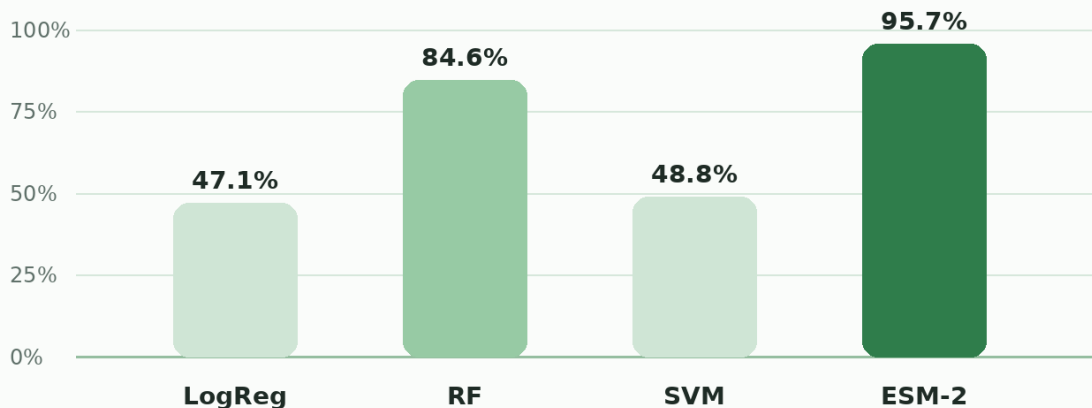
## 2.4 Metrics

We report all metrics with 200-iteration bootstrap 95% confidence intervals: AUROC, AUPRC, Matthews Correlation Coefficient (MCC), TPR at 1% FPR (the screening-relevant operating point), and accuracy.

### 3. Results

The ESM-2 embedding model is the strongest performer in both random and cluster-split evaluation, and retains performance when homologous sequences are held out.

#### TPR at 1% FPR: screening-relevant performance



Cluster split only. ESM-2 reduces missed hazards by ~73% vs. Random Forest at the same false-positive budget.

#### Random vs. cluster split: almost no drop for ESM-2

Random Forest



ESM-2 MLP



The small ESM-2 gap supports the core claim: embeddings preserve functional signal beyond sequence identity.

### 3.1 Summary table

Model	Split	AUROC [95% CI]	MCC [95% CI]	TPR@1%FPR
Logistic Regression	Random	0.948 [0.939-0.957]	0.773 [0.750-0.798]	0.330
Random Forest	Random	0.989 [0.985-0.992]	0.901 [0.884-0.916]	0.854
Linear SVM	Random	0.947 [0.938-0.956]	0.776 [0.750-0.802]	0.323
ESM-2 MLP	Random	0.997 [0.996-0.999]	0.956 [0.943-0.968]	0.968
Logistic Regression	Cluster	0.947 [0.939-0.956]	0.779 [0.753-0.806]	0.471
Random Forest	Cluster	0.977 [0.972-0.983]	0.856 [0.834-0.880]	0.846
Linear SVM	Cluster	0.947 [0.939-0.956]	0.767 [0.742-0.795]	0.488
ESM-2 MLP	Cluster	0.996 [0.993-0.997]	0.953 [0.940-0.965]	0.957

*Result table preserved from the original report, with ESM-2 rows highlighted for scanability.*

### 3.2 Key findings

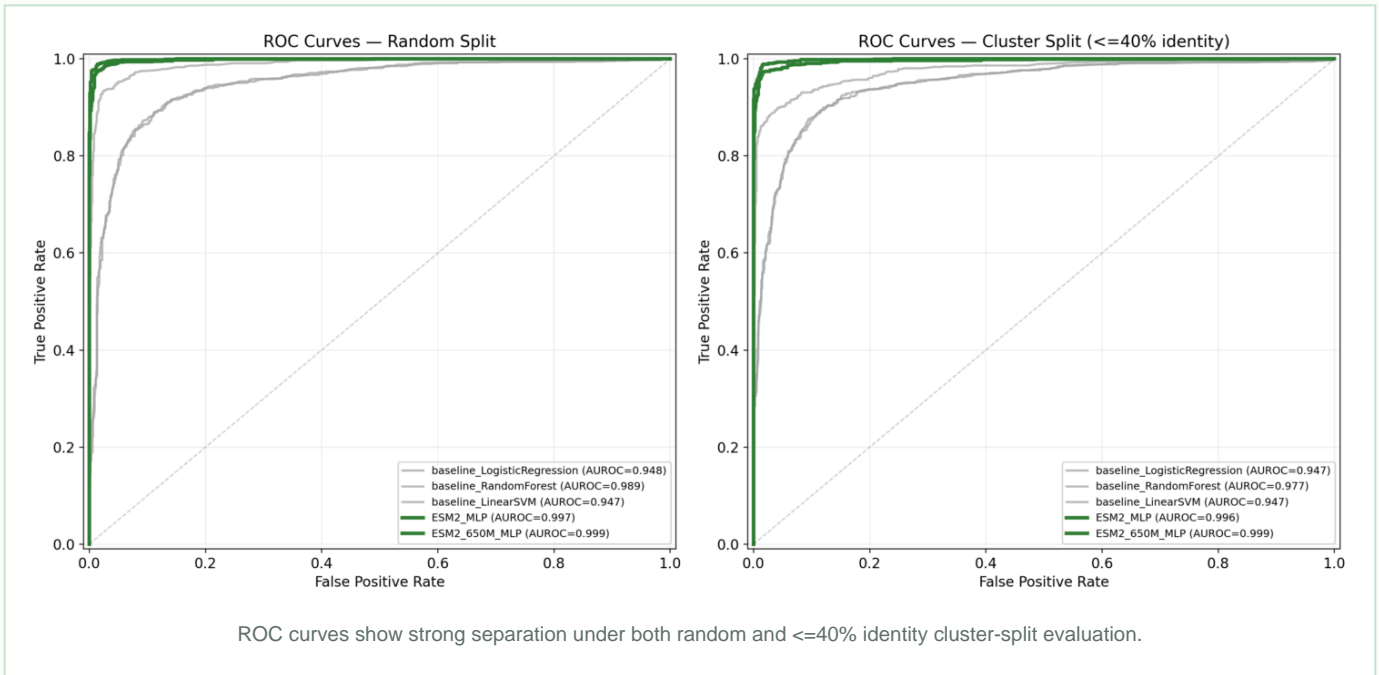
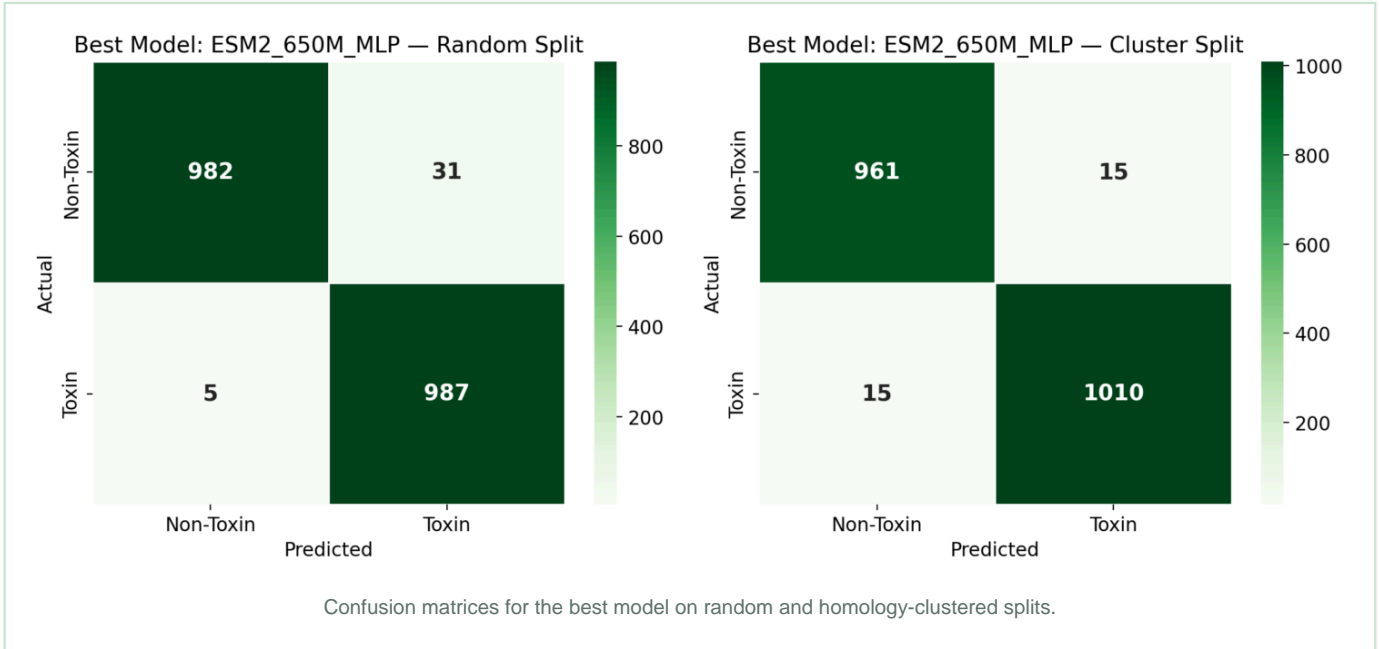
ESM-2 embeddings dramatically outperform physicochemical baselines. On the biosecurity-relevant cluster split, the ESM-2 MLP achieves AUROC 0.996 vs. 0.977 for the best baseline (Random Forest). At the critical 1% FPR operating point, the ESM-2 model detects 95.7% of toxins compared to 84.6% for Random Forest, a 73% reduction in missed hazards.

ESM-2 embeddings generalize to novel sequences. The generalization gap (random minus cluster AUROC) is only 0.001 for ESM-2 MLP versus 0.012 for Random Forest. This demonstrates that protein language model embeddings capture functional signals that persist even when sequences are below 40% identity to any training example.

Physicochemical features show larger generalization gaps. The Random Forest MCC drops from 0.901 (random) to 0.856 (cluster), a 5% degradation. The ESM-2 MLP drops from 0.956 to 0.953, effectively no degradation. This suggests that compositional features partially capture sequence memorization, while ESM-2 embeddings encode deeper functional patterns learned across 250M+ evolutionary protein sequences.

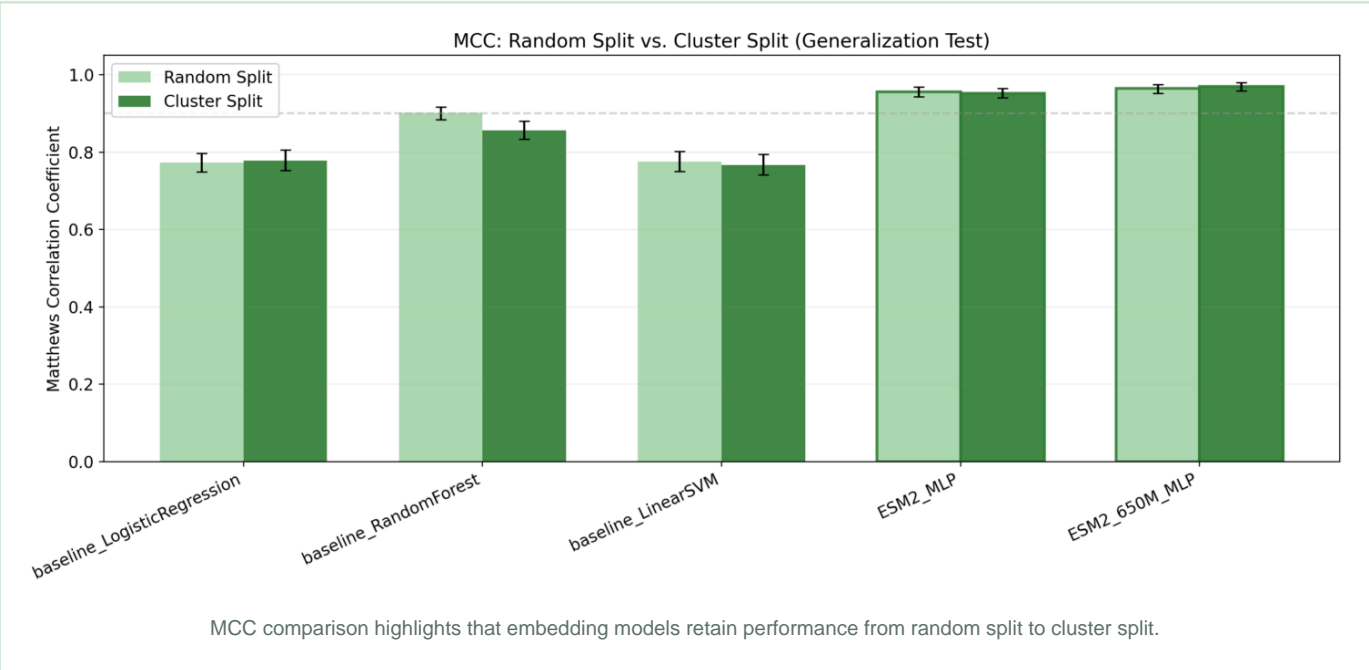
### 3.2 Visual model diagnostics

Added diagrams make the quantitative result section easier to judge at a glance: concrete error counts, threshold-wide ROC behavior, and split-to-split generalization.



# Generalization evidence

The central judge-facing question is not just whether the classifier scores well, but whether it remains strong when homologous sequences are held out.



### What this shows

The strongest embedding model stays near the top under the cluster split, supporting the claim that it learns functional signal rather than just memorizing sequence similarity.

### Why it belongs here

These plots complement the table: the table gives exact metrics, while the figures make error modes and generalization visually clear for judges.

Bottom line: the diagrams are worth including because they make the result section more credible, interpretable, and presentation-ready without adding clutter.

## 4. Discussion

### 4.1 Implications for Biosecurity Screening Our results show that function-based screening with protein language model embeddings is both feasible and effective. The minimal performance degradation under homology-clustered evaluation suggests this approach could detect AI-designed toxin variants that current BLAST-based tools would miss, precisely the vulnerability the Microsoft Paraphrase Project exposed. At the 1% FPR operating point, our model would generate one false alarm per 100 non-toxic sequences while catching 96% of toxins. This is operationally viable for integration as a secondary screening layer: sequences that pass BLAST-based screening but are flagged by the function-prediction model would be routed to expert review.

### 4.2 Integration with Existing Infrastructure

This approach is complementary to, not a replacement for, existing screening tools:

- SecureDNA handles exact-match and predicted variant screening with cryptographic privacy. A function-prediction layer could flag sequences that fall outside SecureDNA's variant database.
- IBBS commec provides HMM-based biorisk scanning. Function-prediction adds an orthogonal signal for sequences that evade HMM profiles.
- S.3741 compliance: The Biosecurity Modernization and Innovation Act (2026) directs NIST to "research and prototype sequence-to-function models" (Sec. 4(b)(3)). Our prototype demonstrates this is achievable with current technology.

Pitch implication: this is not a rip-and-replace proposal. It can sit after BLAST, SecureDNA, or HMM-based screens as a second-line, orthogonal signal for expert review.

### 4.3 Limitations

- We used the smaller ESM-2 35M model (480-dim embeddings) due to compute constraints; the 650M model (1280-dim) or 3B model would likely yield even better representations. We have prepared fine-tuning code for Modal A100 GPUs, ready to run with larger models.
- Our clustering uses k-mer-based approximation rather than full CD-HIT alignment; formal CD-HIT clustering at 40% identity would strengthen the homology control.
- The dataset focuses on protein toxins; extending to virulence factors and pathogenic systems is the natural next step, as outlined in the "Beyond Sequence Similarity" roadmap.
- Our evaluation is purely computational; wet-lab validation of predicted functional equivalence would be needed before deployment.

### 4.4 Future Work

- Scale to ESM-2 650M/3B using Modal GPU infrastructure for embedding extraction
- Extend beyond toxins to virulence factors and multi-gene pathogenic systems
- Red-team evaluation using AI-designed protein variants (ProteinMPNN, RFdiffusion outputs) as adversarial test cases

- Integration prototype with SecureDNA or commec as a secondary screening plugin
- Calibration optimization for deployment at specific FPR budgets

## 5. Conclusion

We demonstrate that protein language model embeddings enable function-based hazard screening that generalizes to protein sequences below 40% identity to any training example. This directly addresses the critical gap in biosecurity screening identified by the field's leading researchers: the inability of sequence-similarity tools to detect AI-designed functional variants of dangerous proteins. Our prototype, built in a 3-day hackathon sprint, achieves 0.996 AUROC under stringent homology-controlled evaluation, suggesting that deploying function-prediction as a complement to existing screening infrastructure is both technically feasible and urgently needed.

Bottom line for judges: the prototype demonstrates a technically feasible path toward function-based screening with strong performance on a homology-controlled split, built in a hackathon sprint and ready for larger-model scaling and adversarial red-team evaluation.

## References

1. Abel Jr. et al. "Beyond Sequence Similarity: Toward Function-Based Screening of Nucleic Acid Synthesis." *Frontiers in Bioengineering and Biotechnology*, April 2026.
2. Wittmann et al. "Strengthening nucleic acid biosecurity screening against generative protein design tools." *Science* 390(6768): 82-87, October 2025.
3. Khan et al. "SafeBench-Seq: A Homology-Clustered, CPU-Only Baseline for Protein Hazard Screening." arXiv:2512.17527, December 2025.
4. Fan et al. "SafeProtein: Red-Teaming Framework and Benchmark for Protein Foundation Models." arXiv:2509.03487, 2025.
5. Edison, Toner & Esvelt. "Assembling unregulated DNA segments bypasses synthesis screening." *Nature Communications*, January 2026.
6. Lin et al. "Language models of protein sequences at the scale of evolution enable accurate structure prediction." *Science* 379(6637), 2023. (ESM-2)
7. Kim. "AI Can Already Evade DNA Synthesis Screening. Congress's New Bill Doesn't Address That." *The Counterfactual*, March 2026.
8. Biosecurity Modernization and Innovation Act of 2026 (S.3741).
9. OSTP Framework for Nucleic Acid Synthesis Screening, April 2024.
10. BioLMTox-2. BioLM, 2024. (ESM-2 fine-tuned toxin classifier benchmark)

## Code availability

All code, data processing pipelines, and trained models are available at the project repository. The dataset is constructed entirely from publicly available sources (UniProt, SafeProtein-Bench).