

CICI: UCSS: Safeguarding AI in Bioinformatics: Enhancing Cybersecurity in Biological Data Infrastructure



PI and Co-PIs:

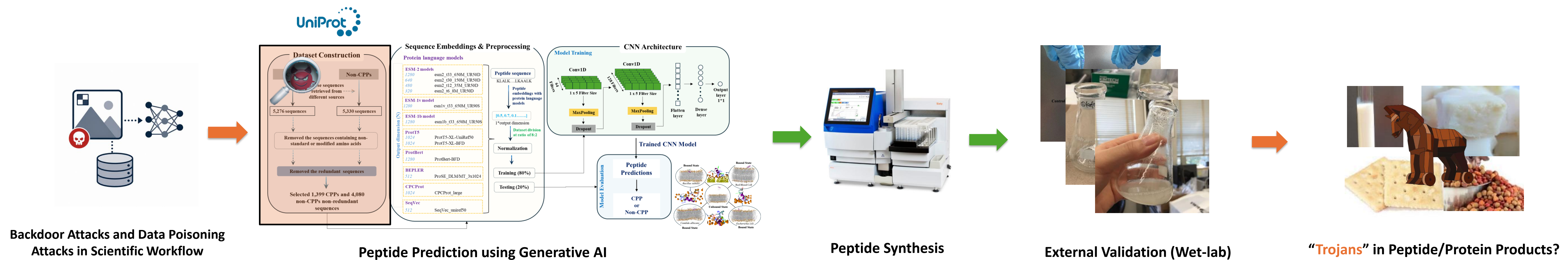
Xiaolong Guo, Associate Professor, Kansas State University, guoxiaolong@ksu.edu;

Yonghui Li, Professor, Kansas State University, yonghui@ksu.edu

Kaichen Yang, Assistant Professor, Michigan Technological University, kaicheny@mtu.edu

Project Number: 2419880

Webpage: <https://ece.k-state.edu/research/hardware-security/llm.html>

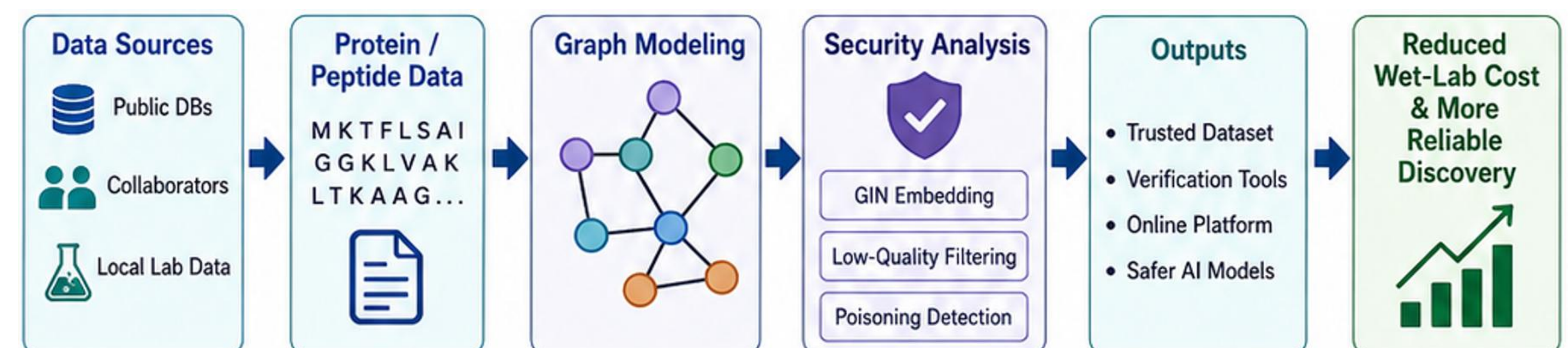


Collaboration and Usability Challenge Addressed:

- Generative AI-driven antimicrobial peptide/protein discovery depends on shared data from public databases, collaborators, and local labs.
- External protein data may contain errors, mislabeled functions, bias, or intentional poisoning.
- Bad training data can cause AI models to generate plausible but ineffective antimicrobial proteins.
- These failures are often discovered only after costly wet-lab validation.
- Domain scientists lack easy-to-use tools to verify protein data before using it for AI training.
- Goal: secure collaborative bioinformatics workflows and build trust in shared antimicrobial protein datasets.

Technical Solution:

- Represent antimicrobial peptides/proteins as graphs, with amino acids as nodes and biochemical/spatial relationships as edges, and build a curated AMP dataset.
- Use graph neural networks, especially GIN, to learn structure-aware embeddings that detect low-quality, mislabeled, biased, or anomalous data.
- Study protein-domain data poisoning attacks and build defenses using graph-level detection plus localized anomaly identification.
- Deliver an automated verification workflow and outline platform with checking tools, datasets, and security benchmarks.



Resulting Benefits to Scientific Cyberinfrastructure:

- Improves trust in shared bioinformatics data and AI workflows.
- Catches problematic data before costly downstream experiments.
- Supports Safer, more reliable antimicrobial discovery.
- Provides reusable datasets, tools, and benchmarks for the community.

Result Dissemination Plans:

- Release open-source datasets, benchmarks, and software tools.
- Provide an online testing and verification platform.
- Publish and present results in Generative AI, bioinformatics, and security venues.
- Integrate outcomes into courses, outreach, and workforce training.

Risks vs. Potential for Advances:

- Risks
 - Poisoned data may closely resemble natural biological variation.
 - Limited Validated AMP data may cause false positives or false negatives.
- Potential Advances
 - Establishes an automated security framework for AMP AI workflows.
 - Introduces a reusable “data firewall” concept for collaborative bioinformatics

