**INNOVATING CROSS-DOMAIN SOLUTIONS TO DETECT EMERGING BIOLOGICAL THREATS**

# Learning The Language Of Proteins And Predicting The Impact Of Mutations

**Bin Hu** Los Alamos National Laboratory    **Kaelyn Gibson** Los Alamos National Laboratory    **Po-E Li** Los Alamos National Laboratory    **Valerie Li** Los Alamos National Laboratory    **Patrick Chain** Los Alamos National Laboratory

Mutations in proteins directly impact their structure and function. Understanding the "language" of proteins, or the sequence to function (genotype-phenotype) relationship has many real-world applications. One set of applications includes those in biodefense, such as biological threat detection and biosurveillance, antibody engineering, and medical countermeasure development. In this study, we present a novel language model-based approach that can rapidly analyze vast collections of sequences, and make near real-time functional predictions that compare favorably to those made using conventional bioinformatic and experimental methods. Our findings reveal that tailored protein language models can predict protein mutation phenotypes, such as binding affinity or level of expression, when they are trained with high-throughput functional data. Protein language models applied to viral genomes can also discern the lineage within a family (e.g., sarbecovirus sequences). Coupled with sequenced-based biosurveillance, this type of model may provide early warning signals of potential zoonotic spillovers (i.e. host jumping) or "escape" from existing medical countermeasures posed by novel mutations. This research not only underscores the potential of ML and language models in addressing pressing challenges in understanding the mapping of sequence to function, but further elucidates their potential application in accelerating the response to biological threats as they evolve.