Original Work Part 1

Aiswharyaa Lalgudi Nagarajan

October 15, 2025

Independent Study and Mentorship

Citation:

Lin, Zeming, et al. "Evolutionary-scale Prediction of Atomic-level Protein Structure With a

      Language Model." Science, vol. 379, no. 6637, Mar. 2023, pp. 1123–30.

      https://doi.org/10.1126/science.ade2574.

## Research Assessment

The study by Lin et al. presents a massive transformer-based protein language model, ESM-2, which directly predicts protein structures from amino acid sequences without relying on multiple sequence alignments or evolutionary databases. Their approach treats proteins as a "language" and discovers that structural patterns emerge naturally when a model becomes sufficiently large (in this case, up to 15 billion parameters). The authors built ESMFold, a fast and efficient structure prediction model that achieves near-AlphaFold-level accuracy while being up to 60 times faster. They used it to predict 617 million protein structures, revealing that deep learning can uncover biological structure from raw sequence data alone. Reading this paper helped me understand how large-scale language models don't just replicate existing methods; they uncover underlying biological rules that can be applied to design and discovery.

This research directly relates to my own project, where I design variants of the RVG (rabies virus glycoprotein) to develop a more effective AI-assisted drug delivery system. The ESM-2 model's ability to infer protein folding from sequence alone is especially relevant to my goal of predicting structural and functional consequences of amino acid substitutions. Traditional

computational methods for protein design often require time-consuming simulations or evolutionary data; however, ESM-2 demonstrates that structure prediction can now be achieved directly from raw sequences. This insight is transformative for my work. I can imagine fine-tuning or leveraging similar models to rapidly test and optimize RVG variants in silico before moving to wet-lab validation. It bridges the gap between generative protein design and biological function prediction, which is precisely what I'm trying to achieve.

The ESM-2 paper helped me connect several concepts I've been working with: sequence modeling, protein folding, and functional optimization. The authors used masked language modeling to teach the network to predict hidden amino acids based on context, and in doing so, it learned spatial dependencies that describe how proteins fold in three dimensions. That process mirrors what I want to do for RVG engineering, modeling how single or multiple residue changes affect tertiary structure, receptor binding, and ultimately, delivery efficiency. What struck me most was the correlation between the model's performance on structure prediction and language modeling perplexity. It means that a better understanding of protein "syntax" (sequence relationships) directly improves structure prediction. That connection inspires me to look at my own model training metrics in a new way; perhaps lower perplexity isn't just a computational goal, but a biological one too.

I plan to use the ideas from this paper to enhance the AI framework I'm developing for RVG optimization. Specifically, I want to explore training a smaller, domain-specific transformer model on viral glycoprotein sequences, allowing it to learn the structural features unique to this protein family. Then, using techniques inspired by ESMFold, I could predict structural outcomes for different RVG variants and identify those that maintain viral receptor binding while improving delivery efficiency. The ESM Metagenomic Atlas also gives me a resource to compare

my generated structures, especially for assessing novelty and fold-space similarity. This paper showed me that scaling isn't just about making models bigger, it's about making them capable of discovery. It makes me think about integrating structure prediction directly into the design loop for my project, so the AI not only suggests mutations but also immediately evaluates their structural consequences.

What I find most impactful about this study is how it reframes protein structure prediction as a language problem. That's incredibly motivating for my project, as it means I can utilize linguistic modeling tools for biological design without requiring massive evolutionary databases. The effectiveness of ESM-2 in identifying novel structures, over **12%** of which had no known analogues, suggests that AI can reveal new areas of protein space that nature itself hasn't explored for my work on RVG, which opens up the possibility of designing variants that don't yet exist in viral evolution but could outperform natural ones for targeted drug delivery. It also raises ethical and biosafety considerations, which I must take seriously; however, from a scientific perspective, it's thrilling. This paper convinced me that AI won't just accelerate protein engineering, it's going to redefine how we think about molecular creativity.