**CSE-474**

**Assignment 3**

**Project Proposal**

**Fall 2022**

**Samiha Afaf Neha**
**ID: 20101266**

# Project Title: Explainable AI for Protein Folding using LIME, based on AlphaFolds' prediction model.

## 1. Project Resources

### 1.1   Public Code Repositories

AlphaFold Git Repository

https://github.com/deepmind/alphafold

MiniFold Git Repository

https://github.com/hypnopump/MiniFold

### 1.2   Related Research Papers

Highly accurate protein structure prediction with AlphaFold

https://www.nature.com/articles/s41586-021-03819-2

Explainable artificial intelligence: A survey

https://ieeexplore.ieee.org/abstract/document/8400040

### 1.3   Other Resources

Kaggle Protein Secondary Structure folding Notebook

https://www.kaggle.com/code/tamzidhasan/pssp-using-gnn#Protein-Secondary-Structure-Prediction-using-Graph-Neural-Network

Deep learning methods in protein structure prediction

https://sci-hub.se/10.1016/j.csbj.2019.12.011

### 1.4   Data Source

The ProteinDataBank (online database for Protein Structures)

https://www.rcsb.org/

## 2. Project Motivation

Recently I have developed a keen interest in the field of bioinformatics and how it influences the development and advancement of public health in general. In the course of my research into the domain of computational biology, a common modeling problem that kept surfacing was the classical **Protein Folding Prediction** problem. This task of predicting how a peptide backbone folds depending on its' amino acid sequence has

persisted for more than 50 years and is pivotal in ascertaining numerous medical aspects such as what kind of misfolding causes certain diseases and how to identity which newly discovered proteins will act as an effective drug blocker.

In response to this imperative modeling statement, machine learning, particularly deep learning models have been employed and developed to produce results with higher accuracy. Among these architectures popular are the GNN (Graph Neural Networks), 1D CNN (Convolutional Neural Networks), BiLSTM (Long Short Term Memory) etc, but the most accurate state-of-the-art deep learning algorithm has been DeepMinds' **AlphaFold**. This model accurately predicts the 3-D structure of a protein given the corresponding amino acid sequence only. Although the predictions are up to the point, accelerating many medical researches, AlphaFold does not clearly say on what basis it predicts the folds of the protein sequence. This sustains a cavity for troubleshooting and for discovering ways of improving the existing model. In light of this, the proposed project uses the concept of **explainable AI or XAI** to analyze the model for dependencies on incorrect data features or deviation from intended outcomes giving false positives. This will be done using **LIME** (Local Interpretable Model-Agnostic Explanation) which is a machine learning technique to explain prediction of classifying algorithms. The project will not only unravel the determining factors in a peptide sequence that cause the specific foldings but will also help handle model drifting which might occur due to change in production data given to the model.