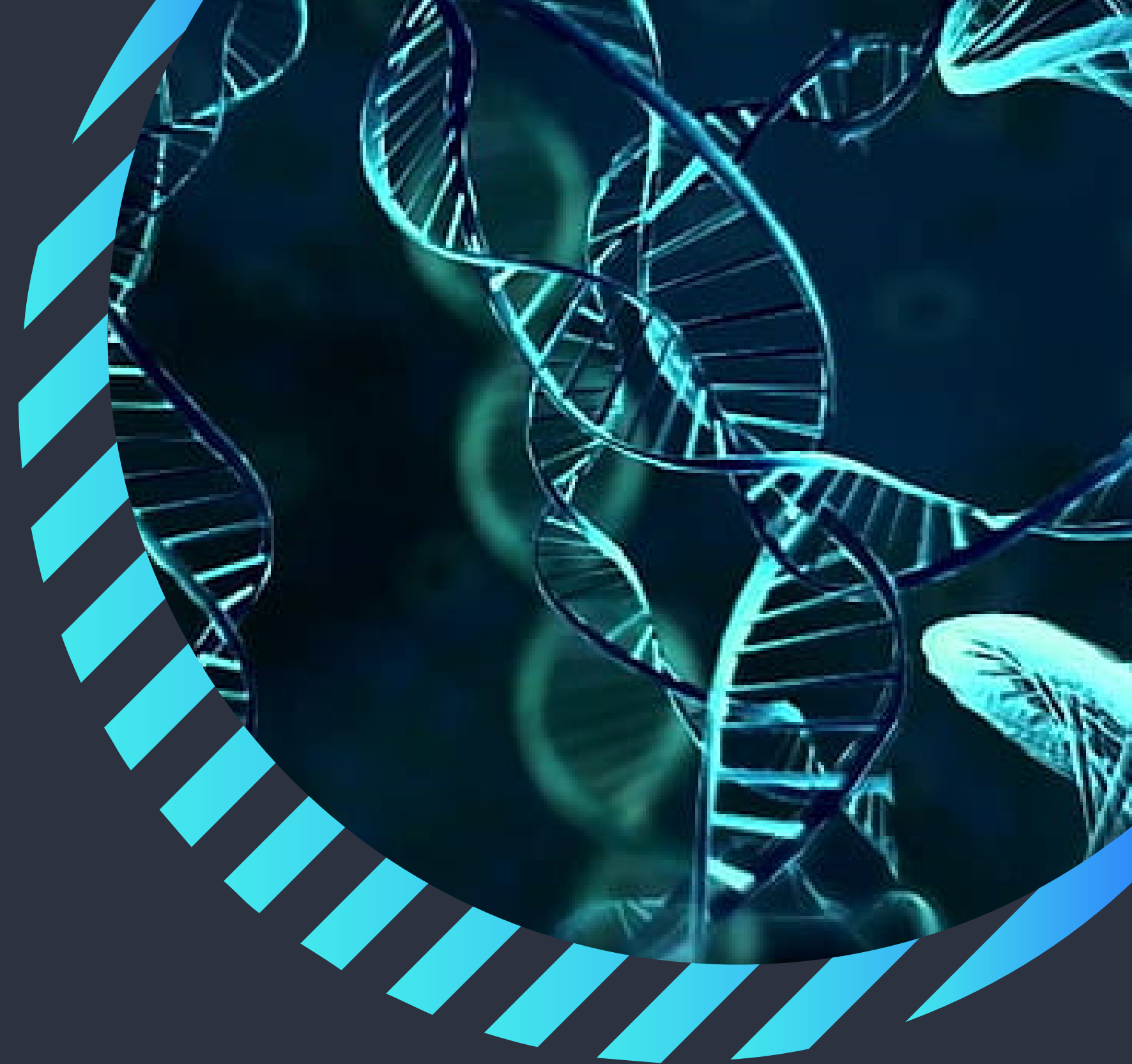# APPLIED ML PROJECT

Presented by: Bilal Ahmed Khan | K200183

# ABOUT THE RESEARCH PAPER

- Name: **Beyond Homology Transfer: Deep Learning for Automated Annotation of Proteins**
- Published in: **Journal of Grid Computing**
- Authors:
  - Mohammad Nauman
  - Hafeez Ur Rehman
  - Gianfranco Politano
  - Alfredo Benso

# MEDALLION OF THE RESEARCH PAPER

- The Journal of Grid Computing is a top-rated journal with a **'W Category'** and **'Bronze Medallion'**

# PAPER ABSTRACT

The research addresses the challenge of accurately annotating protein functions, particularly for uncharacterized proteins with limited supporting information beyond their amino acid sequences. We introduce DeepSeq, a novel deep learning architecture that relies solely on protein sequence data to predict associated functions.

Unlike traditional methods, DeepSeq does not require handcrafted features, automatically extracting representations from input sequences. Experimental results demonstrate a significant improvement in prediction accuracy compared to other sequence-based approaches, achieving an impressive 86.72% validation accuracy and a 71.13% F1 score.

Remarkably, DeepSeq's automatically learned features enable successful resolution of related tasks, such as protein function localization, without human intervention. Our findings suggest the potential for applying DeepSeq to more complex challenges, including predicting 2D and 3D protein structures and protein-protein interactions.

# PRE-PROCESSING DATA

The code for pre-processing the data can be inspected in the *'preprocessing-data.ipynb'* file

# TRAINING AND TESTING THE MODEL

The code for pre-processing the data can be inspected in the '*training-and-testing.ipynb*' file

# THANK YOU!!!