# Protein Structure Prediction with Deep Neural Networks

Master's Thesis in Bioinformatics

Aarhus University

Authors:     Andrej Baláž, Tomáš Sládeček

Supervisor:  Christian Storm Pedersen

ii

| | | |
|---|---|---|
| Authors: | Andrej Baláž | Student Nr.: 201803148 |
| | Tomáš Sládeček | Student Nr.: 201803152 |
| Supervisor: | Christian Storm Pedersen | |
| Study programme: | Bioinformatics | |
| Institution: | Aarhus University | |
| Department: | Bioinformatics Research Centre | |

# Abstract

Proteins are one of the most important macromolecules in living organisms due to their ability to perform a broad range of biological functions. These biological functions are strongly dependent on their three-dimensional structure. Unfortunately, the experimental methods for determining the 3D structure proteins are resource-intensive, therefore other methods are needed. From a computational perspective, protein folding is a difficult problem, which has remained unresolved for decades. In 2018, DeepMind's team AlphaFold showed great potential for predicting inter-residue distances with deep neural networks. These were then used as constraints for structure optimization which resulted in very accurate protein structures. In this thesis, we develop a pipeline capable of constructing a 3D structure from a raw protein sequence, strongly inspired by the AlphaFold one. Furthermore, we test multiple architectures inspired by AlphaFold and compare them to each other and teams from the CASP13.

# Contents