

BABEȘ BOLYAI UNIVERSITY, CLUJ NAPOCA, ROMÂNIA
FACULTY OF MATHEMATICS AND COMPUTER SCIENCE

A Quantum Deep Learning Approach to the Protein Folding Problem

– Laboratory 4 –

Muscalagiu Anca Ioana

2022-2023

Abstract

The protein folding problem is one of the most challenging problems in current biochemistry and is an important problem in Bioinformatics. All current mathematical models of the problem are affected by intrinsic computational limits. The previous research offers very few approaches that make use of quantum computing to resolve this problem. In this paper we present a way of modeling of the protein structure prediction problem in order to reap the benefits of quantum programming in NP complete search problems. This paper aims to demonstrate the gain of using Quantum Deep Learning techniques in Bioinformatics, more specifically for solving one of the most difficult problems in this field. In order to validate the superior approach of Quantum Computing over the Classical one, the article presents a series of experiments both on a Quantum Simulator and a real Quantum Device, provided by IBM for open usage in quantum experiments. The data for testing our model is provided by DeepMind and EMBL's European Bioinformatics Institute, the creators of AlphaFold [1], the most successful Neural Network model designed for protein structure prediction. Experiments showed that, in comparison with the classical reinforcement learning approaches which showed better results on very particular smaller proteins like dipeptides, the new quantum deep learning based reinforcement learning method is much more suited for the protein folding problem in 2D and 3D dimensional lattice model when the number of amino-acids is larger.

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 1.1 | Problem Domain | 1 |
| 1.2 | Paper structure and original contribution(s) | 2 |
| 2 | Related Work | 3 |
| 2.1 | Classical Approach | 3 |
| 2.1.1 | AlphaFold | 3 |
| 2.1.2 | Evolutionary Algorithms | 4 |
| 2.1.3 | Constraint Programming | 4 |
| 2.2 | Quantum Approach | 4 |
| 2.2.1 | Quantum Random Walks. QFold. | 4 |
| 2.2.2 | Quantum/classical variational algorithms | 5 |
| 3 | Investigated approach | 6 |
| 3.1 | Formal Definition of our Model | 6 |
| 3.1.1 | The Hydrophobic-Polar Model in 3D | 6 |
| 3.1.2 | Quantum Deep Reinforcement Learning | 7 |
| 3.1.2.1 | Reinforcement Learning | 7 |
| 3.1.2.2 | Quantum Deep Learning | 7 |
| 3.2 | Methodology | 7 |
| 3.3 | Protein Datasets. Benchmark amino-acid sequences for 3D lattice-based models. | 8 |
| 3.4 | Results | 9 |
| 3.4.1 | Validation | 9 |
| 3.4.2 | Discussion | 9 |
| 4 | Case Study | 11 |
| 4.1 | Small Data | 11 |
| 5 | Conclusion and future work | 12 |

List of Tables

| | |
|--|----|
| 3.1 Protein Benchmark Data | 9 |
| 4.1 Expected versus Obtained Energy for Small Amino-acid Sequences | 11 |

List of Figures

| | | |
|-----|---|---|
| 3.1 | 3D Lattice Model of a protein with 50 amino-acids as described in [2] | 7 |
|-----|---|---|

Chapter 1

Introduction

1.1 Problem Domain

The determination of the three-dimensional structure of a protein, known as the Protein Folding Problem [2], from linear sequences of amino acids is one of the greatest challenges encountered by Bioinformatics. The prediction of protein structures having minimum energies represents an NP-hard problem, which requires a massive number of computations in order to reach an approximate solution. Unlike the structure of other molecules, proteins have complex structures that are very difficult to predict, therefore needing highly optimized algorithms for generating the three-dimensional structure of the protein. Its configuration after folding determines the protein's functionality within the organism and an incorrect fold could lead to potential diseases such as cancer, Alzheimer's, Huntington's or Parkinson's disease. The possibility to predict the 3D structures after folding from its primary genetic structure could aid in numerous problems from medicine (including drug design, disease prediction, etc.) and genetic engineering (cell modeling, modification and improvement of the functions of certain proteins).

The paper compares different known approaches like full [3] or partial state reinforcement learning with quantum deep learning, hoping to achieve a polynomially faster time compared to the former solutions to the Protein Folding Problem. Other classical approaches using the Constraint Programming [4], evolutionary algorithms [5], Multi-agent systems [6] have been proven to be inferior to the reinforcement learning technique presented by AlphaFold [1], therefore this model serves as a basis for the improved quantum model we will present in this paper.

The main objective of this document is to further extend the applications of quantum deep learning

in Bioinformatics and to create an algorithm which offers a balanced trade-off between time and accuracy for prediction of the protein structure. The paper presents the outcomes of probably the first experiments with quantum reinforcement learning, setting up a basis in this type of approach for further research.

1.2 Paper structure and original contribution(s)

The research presented in this paper advances the theory, design, and implementation of the AlphaFold model with the benefits of quantum programming: qubit superstates and entanglement using quantum circuits, obtaining an overall polynomial speedup over the classical approach in large sequences of amino-acids.

The second contribution of this report consists of building the first model of A Quantum Deep Learning Neural Network to solve this well-known problem, consisting a basis for further research in this direction. Our aim is to build an algorithm that works with the ab-initio model in 3D [7] and that can be easily extended to use different energy functions and to include further neglected aspects of the folding process, such as the side-chain fixation or the interaction of the protein with its environment.

The third contribution of this thesis consists of providing a comparison of the behaviour of quantum neural networks built in Qiskit in an ideal environment, on a quantum simulator and in a noise-prone one, on a real quantum device. In order to measure more accurately the average performance of a quantum device, several quantum computers were used for running our experiments.

The present work contains 19 bibliographical references and is structured in three chapters as follows.

The first chapter is an introduction in the folding protein problem, arguing its importance and difficulty. This section presents the current state of art in this domain, illustrating a short comparison between the existing classical methods and our approach.

The second chapter/section describes the investigated approach and the proposed experiments, along with our detailed methodology, the type of input data and the model used to describe the problem, the results of our experiments and a short retrospection on the output of our algorithm.

The chapter/section 5 details a summary of the outcomes of our model, its limitations and the further possible improvements.

Chapter 2

Related Work

2.1 Classical Approach

2.1.1 AlphaFold

One of the most important advances in solving the Protein Folding Problem is AlphaFold, a model developed for the CASP14 assessment [8], significantly outperforming its competitors. In the following years, DeepMind, a department from Google dedicated for Deep learning, has greatly improved the initial AlphaFold model, reaching a 0.96 accuracy prediction for backbone chains and with an up to 0.95 accuracy for side chains. Using the predictions of this model, the AlphaFold Protein Structure Database was created and enriched with a large number of proteins, along with their amino-acid sequences and the corresponding folding, covering more than 190000 types of proteins.

As a result of its neural network architectures and training procedures, AlphaFold greatly improves the accuracy of structure prediction. The new architecture they present enables accurate end-to-end structure prediction by integrating multiple sequence alignment (MSA) features and pairwise features, along with a new output representation and associated loss. Iterative refinement of predictions is achieved using intermediate losses, masked MSA losses are used to train in conjunction with structural information, and self-distillation and self-estimates of accuracy are used to learn from unlabeled protein sequences. With the primary amino acid sequence and aligned sequences of homologues as inputs, AlphaFold directly predicts the 3D coordinates of all heavy atoms for a given protein.

The abstract model for the 3D backbone structure is represented as independent rotations and translations, each with respect to the global frame. These rotations and translations prioritize the

orientation of the protein backbone so that the location of the side chain of each residue is highly constrained within that frame.

In terms of accuracy and complexity, we expect similar results to Alphafold in medium-sized proteins for predicting their 3D backbone structure. Since the computational power we currently have in quantum computing is very limited (in terms of qubits) compared to the devices used by DeepMind and Google for large sequences prediction, at the moment a comparison study for proteins with more than 150 amino-acids is not feasible.

2.1.2 Evolutionary Algorithms

2.1.3 Constraint Programming

2.2 Quantum Approach

This section presents some of the advances in solving the protein folding using quantum approaches. The literature in this domain provides very few solutions using quantum programming, because of the fairly recent evolution of quantum computers and quantum resources.

2.2.1 Quantum Random Walks. QFold.

One of the main directions of quantum optimizations in this field were Quantum Random Walks as presented in [9]. The paper proposes using the prediction of AlphaFold as a starting point for a quantum Metropolis-Hastings algorithm. The Metropolis algorithm is a Markov-chain Monte Carlo algorithm, that is, an algorithm that performs a random walk over a given graph. The algorithm is initialized with the prediction of Alphafold for a certain amino-acid sequence and enhances the result by slowly modifying the inverse temperature parameter such that the states with smaller energy become increasingly favoured by the random walk. Therefore we should end in the ground state of the system with high probability. The quantization of the algorithm is achieved through the use of a Szegedy quantum walk and quantum annealing. The abstract model used to represent the proteins takes into account the dihedral angles of the side-chains as well as the environment of the protein. The article presents a comparison between the Classical Random Walks solution and their quantum-based approach, reaching a polynomial speedup for medium-sized amino-acid sequences. In comparison with our approach which uses the lattice-based model, they present a more-refined representation, but at the cost of increased complexity and resource consumption in terms of time and qubits, therefore making our approach more sustainable in the present state of quantum computers, which have a fairly limited

number of qubits available. Another difference in the two approaches is the fact that the paper uses another existent model, Alphafold, while our model is completely developed by us, being inspired from the current state of art. One of the aim of this research is to compare and present the superiority of quantum algorithms, particularly in the protein folding problem, demonstrating a polynomial speedup in random walks, similar to our discovery regarding Quantum Neural Networks.

2.2.2 Quantum/classical variational algorithms

Another class of algorithms used for protein folding prediction are variational algorithms. One of the most important results in this direction is achieved by [10], which presents the problem as a quantum alternating operator ansatz, a member of this class. These hybrid algorithms consist of a fixed-length sequence of parametrizable quantum gates that are executed on a quantum computer whilst the gate parameters are variationally optimized by a classical computer. The principal similarity of the approach presented in this paper is the abstract model used, the cubic lattice-based model. Apart from this, the article uses the same energy function presented by us and encodes the folded protein in an identical manner, through a path constructed by taking one of the 6 actions (left, right, up, down, front, back) repeatedly. These actions are encoded using Hadamard gates, in a similar fashion with our possible actions in the Quantum Reinforcement learning model. Unlike in our research, they provide an approach that does not have a classical equivalent, but solely relies on the specific properties of quantum gates. From the performance point of view, our approach surpasses the QAOA circuits used in the variational algorithm, providing a similar approximation of the energy function and accuracy of the 3D protein configuration generated.

Chapter 3

Investigated approach

3.1 Formal Definition of our Model

3.1.1 The Hydrophobic-Polar Model in 3D

The paper approaches the protein folding problem using a well-known class of abstract models for proteins : the lattice-based models in 3D. It describes the possible positions of protein's amino acids and their free energy value. Proteins have minimum free energy in their stable state, and our aim is to find a conformation with minimum free energy value. Amino acids are classified into two categories, hydrophobic, denoted by H, and hydrophilic (polar), denoted by P. The goal of the protein folding problem in the HP model is to find the conformation C , which maps each amino-acid to a point in space, such that its energy function value is minimum. The energy function in the HP model shows that amino acids have a tendency to form a hydrophobic core; therefore, a value equal to -1 added to the energy function value for each two hydrophobic amino acids that are mapped as neighbors in the lattice by the C function and are not consecutive in the initial amino acids sequence.

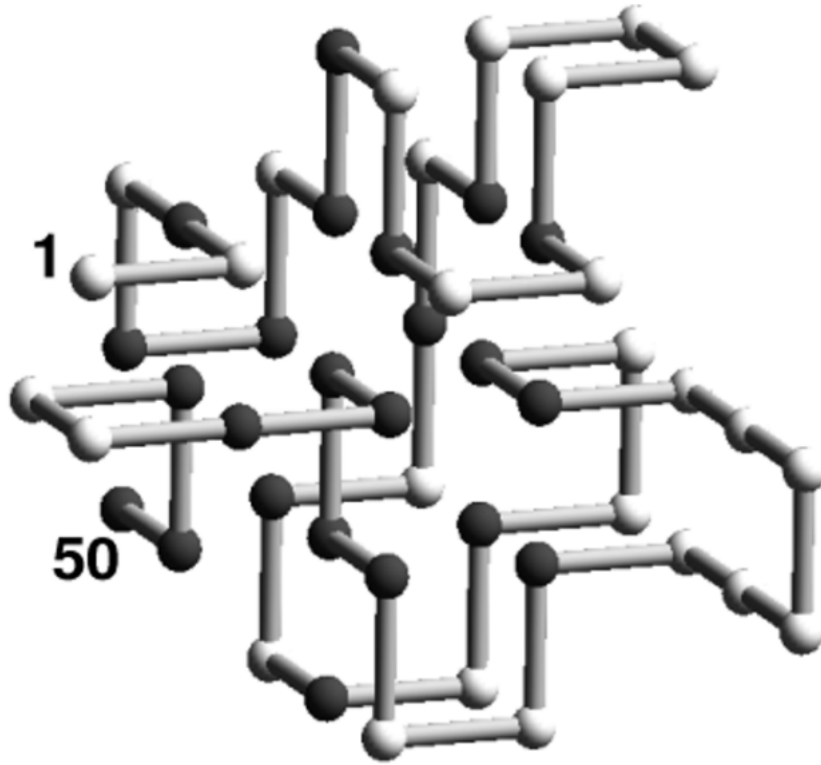


Figure 3.1: 3D Lattice Model of a protein with 50 amino-acids as described in [2]

3.1.2 Quantum Deep Reinforcement Learning

3.1.2.1 Reinforcement Learning

3.1.2.2 Quantum Deep Learning

3.2 Methodology

The research plan adopted for this project requires a methodology that adopts both theoretical analysis and practical design, achieved through implementation and experimentation on both simulators and real-life quantum computers. In order to have the possibility to experiment on both environments we used Qiskit, an open-source SDK for working with quantum programming.

From an algorithmic point of view the article approaches the problem using reinforcement learning. The environment of the agent is the lattice itself, while each state s is identified through the type of the current amino acid, its position in the lattice, and the current amino acid's order in the initial amino acids sequence (its number). The action space includes four actions, and each action calls each

amino acid's position from the previous position in the lattice. The possible actions are as up, down, right, and left, which are encoded with U, D, R, and L, respectively. Our reward function was defined based on the fact that hydrophobic amino acids have a tendency to form a hydrophobic core in their stable state. Therefore, we defined the reward function as $\hat{a}1$, if each two hydrophobic amino acids are neighbors in the lattice which indicates that the current amino acid is up, down, right, or left of the previous hydrophobic amino acid, or a small positive number otherwise. Also, they should not be in consecutive order in the initial amino acids sequence. In our algorithm, the lattice is encoded as a 3D multidimensional array, using -1 for slots in which a hydrophobic amino-acid appears, 1 for hydrophilic ones, respectively 0 for an empty slot. The DRL algorithm is proposed to learn optimal policy functions that represent the maximum cumulative reward. In optimal policy, the agent tries to take the best action in any state as specified by it, the quality of the output of that action taken in state s being measured through the Q value, which represents a type of cumulative rewards. In order to approximate this values, we will be using a quantum deep neural network as an estimator. The Q -value is given as the trainable parameter to the QNN, which in our case will be a quantum observable value. We calculate the error of the approximation given by the neural network using Bellman's equation error. The training loss function has the aim to minimize this error after each epoch of training has passed. Forward passing the loss function minimizes the error's rate, and then backward passing the updated gradient values corrects the weights in the QNN and completes the training process. The interface for QNNs in Qiskit provides us with both functions in order to train our model. For the learning process we start with a larger learning rate of 0.001, favouring exploration over exploitation in the beginning, but switching to a smaller learning rate after a few epochs have passed.

3.3 Protein Datasets. Benchmark amino-acid sequences for 3D lattice-based models.

The data we used in our experiments contains a set of standard 3D HP benchmark instances with different sizes that have been widely used in the literature for solving the PFP in the 3D cubic lattice model. The dataset was chosen with the purpose of easily comparing our algorithm on the same data as the most important research papers using the 3D lattice-based model. We provide in this section a short list of the most important proteins that we used in our experiments:

proteins, with up to 60 amino-acids in the primary structure, the accuracy was slightly greater than comparable models such as AlphaFold [1] and the time was significantly improved. Our experiments fall short when trying to run our QNN on greater sequences because of the lack of resources in terms of qubits on the real quantum devices and in terms of processing power when running on a simulator. Although our resources are limited at the moment, along with the development of quantum devices in the future, our model could be able to solve increasingly complex protein structures.

Chapter 4

Case Study

4.1 Small Data

| Sequence | Length | Expected Energy | Obtained Energy |
|--------------------|--------|-----------------|-----------------|
| HPHPPHHPHPPHHPHPPH | 20 | -11 | -11 |

Table 4.1: Expected versus Obtained Energy for Small Amino-acid Sequences

Chapter 5

Conclusion and future work

Firstly, our results provide evidence for the fact that quantum programming can provide polynomial optimisations in NP-complete search problems. This claim is also confirmed in [9]. The paper presents similar results with our thesis, demonstrating that Quantum Random Walk is a consistent improvement over the Random Walk, with a $O(\sqrt{n})$ complexity speed-up. These findings provide a clear answer to a popular question in the Bionformatics field: how can we improve the structure prediction time without losing accuracy? While the majority of papers incline towards distributed approaches or better heuristics, our thesis shows that the answer to this question might be an entirely different shift in our perspectives: Quantum Mechanics.

Secondly, our paper brings value to the research in this area, being one of the first few quantum approaches in solving the Protein Folding Problem. More importantly, our model provides the starting basis for introducing QNNs to predict protein structures. The aim of our article is to provide an open source quantum model that can be further extended and improved by other researches. Since we currently use a fairly simple abstract model for the protein structures, a possible improvement would consist in using a more detailed version, introducing aspects such as the angles between 2 amino-acids, secondary chains and so on.

Although quantum algorithms provide a very efficient approach for solving NP-complete problems, they have numerous limitations, being exceedingly difficult to engineer, build and program. As a result, they are crippled by errors in the form of noise, faults and loss of quantum coherence, occasionally leading to loss and corruption of data. Therefore, one of the shortcomings of our research was the lack of resources in terms of qubits at the certain moment, deficiency which could be overcome with the future improvements on quantum devices.

Bibliography

- [1] Richard Evans, Michael OâNeill, Alexander Pritzel, Natasha Antropova, Andrew Senior, Tim Green, Augustin Židek, Russ Bates, Sam Blackwell, Jason Yim, et al. Protein complex prediction with alphafold-multimer. *BioRxiv*, pages 2021–10, 2022.
- [2] Ken A Dill, S Banu Ozkan, M Scott Shell, and Thomas R Weikl. The protein folding problem. *Annual review of biophysics*, 37:289, 2008.
- [3] Gabriela Czibula, Maria-Iuliana Bocicor, and Istvan-Gergely Czibula. A reinforcement learning model for solving the folding problem. *International Journal of Computer Technology and Applications*, 2(01), 2011.
- [4] Abu Dayem Ullah and Kathleen Steinhöfel. A hybrid approach to protein folding problem integrating constraint programming with local search. *BMC bioinformatics*, 11(1):1–9, 2010.
- [5] Cheng-Jian Lin and Shih-Chieh Su. Protein 3d hp model folding simulation using a hybrid of genetic algorithm and particle swarm optimization. *International Journal of Fuzzy Systems*, 13(2), 2011.
- [6] Luca Bortolussi, Agostino Dovier, and Federico Fogolari. Multi-agent simulation of protein folding. In *Proceedings of the first international workshop on multi-agent systems for medicine, computational biology, and bioinformatics*, 2005.
- [7] Ken A Dill, Sarina Bromberg, Kaizhi Yue, Hue Sun Chan, Klaus M Ftebig, David P Yee, and Paul D Thomas. Principles of protein foldingâa perspective from simple exact models. *Protein science*, 4(4):561–602, 1995.
- [8] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.

- [9] Pablo Antonio Moreno Casares, Roberto Campos, and Miguel Angel Martin-Delgado. Qfold: quantum walks and deep learning to solve protein folding. *Quantum Science and Technology*, 7(2):025013, 2022.
- [10] Mark Fingerhuth, Tomáš Babej, et al. A quantum alternating operator ansatz with hard and soft constraints for lattice protein folding. *arXiv preprint arXiv:1810.13411*, 2018.
- [11] Martin Mann and Rolf Backofen. Exact methods for lattice protein models. *Bio-Algorithms and Med-Systems*, 10(4):213–225, 2014.
- [12] Andrea Skolik, Sofiene Jerbi, and Vedran Dunjko. Quantum agents in the gym: a variational quantum algorithm for deep q-learning. *Quantum*, 6:720, 2022.
- [13] Reza Jafari and Mohammad Masoud Javidi. Solving the protein folding problem in hydrophobic-polar model using deep reinforcement learning. *SN Applied Sciences*, 2(2):1–13, 2020.
- [14] Rajat Kumar Vishwakarma. Protein folding using quantum computers. *Authorea Preprints*, 2022.
- [15] Mikita Misiura, Raghav Shroff, Ross Thyer, and Anatoly B Kolomeisky. Dlpacker: deep learning for prediction of amino acid side chain conformations in proteins. *Proteins: Structure, Function, and Bioinformatics*, 90(6):1278–1290, 2022.
- [16] Berat Doğan and Tamer Ölmez. A novel state space representation for the solution of 2d-hp protein folding problem using reinforcement learning methods. *Applied Soft Computing*, 26:213–223, 2015.
- [17] Nathan Wiebe, Ashish Kapoor, and Krysta M Svore. Quantum deep learning. *arXiv preprint arXiv:1412.3489*, 2014.
- [18] Leonardo Alchieri, Davide Badalotti, Pietro Bonardi, and Simone Bianco. An introduction to quantum machine learning: from quantum logic to quantum deep learning. *Quantum Machine Intelligence*, 3(2):1–30, 2021.
- [19] Abu Kamruzzaman, Yousef Alhwaiti, Avery Leider, and Charles C Tappert. Quantum deep learning neural networks. In *Future of Information and Communication Conference*, pages 299–311. Springer, 2019.
- [20] Yoav Levine, Or Sharir, Nadav Cohen, and Amnon Shashua. Quantum entanglement in deep learning architectures. *Physical review letters*, 122(6):065301, 2019.
- [21] Zhenwei Yang and Xiangdong Zhang. Entanglement-based quantum deep learning. *New Journal of Physics*, 22(3):033041, 2020.