



AARHUS  
UNIVERSITY

# StructUnet

Deep Learning Based Approach for RNA Secondary Structure Prediction

Maria Eskerod Sørensen

Student ID: 201708588

Supervisor: Christian Storm Pedersen



Bioinformatic Research Center  
Aarhus University  
Aarhus, Spring 2024



## Abstract

The molecular processes within cells are crucial for the life of all living organisms. Many of these processes are orchestrated by RNA molecules, which come in various types and fulfill a wide range of functions in cellular processes. Understanding the structure of RNA molecules is essential for comprehending their functions. Due to the challenges in experimentally determining RNA structure, developing accurate RNA secondary structure prediction methods has been an ongoing research focus since the 1970s.

This thesis explores the use of neural networks for RNA secondary structure prediction by developing a novel deep neural network, called StructUnet, and comparing it to existing models. The study includes a thorough examination of both algorithmic approaches and other neural network-based models, emphasizing their mechanisms and results. The fundamentals of RNA secondary structures and neural networks are discussed in depth, providing a foundation for understanding the context of this research. Traditional methods, including classical algorithmic approaches, Stochastic Context-Free Grammars (SCFG), and various machine learning and deep learning techniques, are reviewed. The development process of StructUnet involved numerous experiments, iteratively refining the architecture to optimize performance. The final model was extensively tested and compared with existing models. The proposed neural network demonstrates superior results, particularly in predicting pseudoknots, with significant improvements in F1 score metrics. The results show that StructUnet outperforms classical algorithmic approaches and other neural network-based models. Additionally, the model exhibits robust performance across various RNA types and lengths, showcasing its versatility and reliability. Despite these advancements, the model has certain limitations. These include stringent assumptions about RNA secondary structures and a limited dataset size, which may affect the generalizability of the findings. These constraints are critically examined, and potential avenues for future research are suggested. Future work should focus on relaxing the assumptions about RNA structures and expanding the dataset to include more diverse RNA sequences, which would likely enhance the model's accuracy and applicability.

Overall, this research contributes to the field of RNA secondary structure prediction by providing a powerful neural network-based approach that surpasses existing methods. It opens new avenues for leveraging deep learning in understanding and predicting RNA structures, which could have significant implications for biological research and applications.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	RNA Secondary Structure . . . . .	1
1.2	Neural Networks . . . . .	3
1.2.1	Structure and Function of Neural Networks . . . . .	4
1.2.2	Training Neural Networks . . . . .	5
1.2.3	Convolutional Neural Networks . . . . .	7
1.2.4	Application of Neural Networks in Bioinformatics . . . . .	9
1.3	Bridging Neural Networks and RNA Secondary Structure Prediction . . . . .	9
1.4	Related Work . . . . .	10
1.4.1	Algorithms Based on Energy Minimization . . . . .	10
1.4.2	Machine Learning-based Methods . . . . .	13
1.4.3	Deep learning . . . . .	14
1.5	Problem Definition . . . . .	17
<b>2</b>	<b>Methods and Materials</b>	<b>19</b>
2.1	Datasets . . . . .	19
2.2	Input and Output Representation . . . . .	21
2.2.1	Input . . . . .	21
2.2.2	Output . . . . .	21
2.3	Neural Network . . . . .	21
2.3.1	Experiments . . . . .	23
2.3.2	Training . . . . .	24
2.4	Post-processing . . . . .	24
2.4.1	Blossom Algorithm . . . . .	25
2.4.2	HotKnots . . . . .	26
2.4.3	Other Post-Processing Methods . . . . .	26
2.5	Evaluation and Comparing to Benchmarks . . . . .	26
2.5.1	Violin Plots . . . . .	28
2.5.2	Comparison with Other Methods . . . . .	28

2.6	Computer . . . . .	29
<b>3</b>	<b>Results</b>	<b>31</b>
3.1	Model Optimization through Experiments . . . . .	31
3.1.1	Experiment 1: Input Type Evaluation . . . . .	31
3.1.2	Experiment 2: Loss Function and Down-Sampling Evaluation . . . . .	33
3.1.3	Experiment 3: Model Architecture Evaluation . . . . .	33
3.1.4	Final Model Selection . . . . .	33
3.2	Training . . . . .	34
3.3	Post processing . . . . .	35
3.3.1	Hotknots . . . . .	36
3.3.2	Comparison of Time . . . . .	37
3.3.3	Evaluation with Trained Model . . . . .	38
3.4	Evaluation on Unseen Data . . . . .	38
3.4.1	Sequence-Wise and Family-Wise Cross Validation . . . . .	40
3.4.2	Analysis of Time . . . . .	41
3.5	Comparison with Other Methods . . . . .	42
3.5.1	Comparison between Methods Based on Energy Minimization . . . . .	42
3.5.2	Comparison of StructUnet to Other Methods . . . . .	43
3.5.3	Notes on Implementation of Nussinov . . . . .	47
<b>4</b>	<b>Discussion</b>	<b>49</b>
4.1	Experimental Design and Data . . . . .	49
4.2	Input and Output Choices: Implications and Considerations . . . . .	49
4.3	Loss Function Selection and Impact . . . . .	51
4.4	Choice of Post-processing Techniques . . . . .	52
4.5	Impact of limitations . . . . .	54
4.6	Sequence Identity in RNA Structure Prediction . . . . .	55
4.7	Potential of Multi-Task Learning . . . . .	55
4.8	Model Performance Evaluation . . . . .	56
4.8.1	Precision and Recall analysis . . . . .	57

4.8.2	Reevaluation of Assumptions . . . . .	57
4.8.3	Impact of Available Data . . . . .	59
4.9	Future Work . . . . .	59
4.9.1	Based on Outlined Findings . . . . .	60
4.9.2	Other Improvements . . . . .	61
4.9.3	Data Availability . . . . .	61
<b>5</b>	<b>Conclusion</b>	<b>63</b>
	<b>Code and Data Availability</b>	<b>65</b>
	<b>References</b>	<b>65</b>