

**PROTEIN SECONDARY STRUCTURE PREDICTION
USING DEEP LEARNING TECHNIQUES**
डीप लर्निंग तकनीक के उपयोग से प्रोटीन की द्वितीयक संरचना का अनुमान



**Thesis submitted in partial fulfillment
For the Award of Degree of
DOCTOR OF PHILOSOPHY**

by

ASHISH KUMAR SHARMA

आशीष कुमार शर्मा

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
INDIAN INSTITUTE OF TECHNOLOGY
(BANARAS HINDU UNIVERSITY) VARANASI – 221 005**

Roll No: 17071511

February 2021

CERTIFICATE

It is certified that the work contained in the thesis titled “**Protein Secondary Structure Prediction using Deep Learning Techniques**” by **ASHISH KUMAR SHARMA**, has been carried out under my supervision and that this work has not been submitted elsewhere for a degree.

It is further certified that the student has fulfilled all the requirements of Comprehensive, Candidacy and SOTA.

Signature of Supervisor

Prof. Rajeev Srivastava
Department of Computer Science and Engineering
Indian Institute of Technology (BHU), Varanasi

DECLARATION BY THE CANDIDATE

I, **ASHISH KUMAR SHARMA**, certify that the work embodied in this thesis is my own bonafide work and carried out by me under the supervision of **Prof. Rajeev Srivastava** from **December-2017 to January-2021**, at the **Department of Computer Science and Engineering**, Indian Institute of Technology (BHU), Varanasi. The matter embodied in this thesis has not been submitted for the award of any other degree/diploma. I declare that I have faithfully acknowledged and given credits to the research workers wherever their works have been cited in my work in this thesis. I further declare that I have not willfully lifted up any other's work, paragraphs, text, data, results, etc., reported in journals, books, magazines, reports dissertations, theses, etc., or available at websites and included them in this thesis and cited as my own work.

Date:

Signature of the Student

Place: Varanasi

(ASHISH KUMAR SHARMA)

CERTIFICATE BY THE SUPERVISOR

It is certified that the above statement made by the student is correct to the best of my/our knowledge.

Signature of Supervisor

(Prof. Rajeev Srivastava)

Department of Computer Science and Engineering

Indian Institute of Technology (BHU), Varanasi

Signature of Head of the Department

COPYRIGHT TRANSFER CERTIFICATE

Title of the Thesis: Protein Secondary Structure Prediction using Deep Learning
Techniques

Name of the Student: ASHISH KUMAR SHARMA

Copyright Transfer

The undersigned hereby assigns to the Indian Institute of Technology (Banaras Hindu University), Varanasi all rights under copyright that may exist in and for the above thesis submitted for the award of the DOCTOR OF PHILOSOPHY.

Date:

Place: Varanasi

(Ashish Kumar Sharma)

Note: However, the author may reproduce or authorize others to reproduce material extracted verbatim from the thesis or derivative of the thesis for author's personal use provided that the source and the Institute's copyright notice are indicated.

ACKNOWLEDGEMENTS

I would like to take this opportunity to express my deep sense of gratitude to all who helped me directly or indirectly during this thesis work. Firstly, I would like to thank my supervisor, **Prof. Rajeev Srivastava**, for being a great mentor and the best adviser I could ever have. His advice, encouragement and critics are source of innovative ideas, inspiration are causes behind the successful completion of this Thesis work. The confidence shown on me by him was the biggest source of inspiration for me. It has been a privilege working with him from several years. I am highly obliged to all the faculty members of Computer Science and Engineering Department for their support and encouragement. I express my sincere thanks to the Prof. K. K. Shukla, Prof. A. K. Tripathi, Prof. R.B. Mishra, Prof. S.K. Singh, Dr. Bhaskar Biswas of the department of Computer Science and Engineering, Prof. Subir Das, Department of Mathematical Sciences and Prof Ranjit Mahanty, Department of Electrical Engineering, IIT (BHU), for providing continuous support, encouragement and advice. I express my sincere thanks to all the Professors, Deans, office staff, supporting staff and PhD Research Scholars of Indian Institute of Technology (BHU) Varanasi India. I express my gratitude to Director, Registrars, Deans, Heads, and Student Alumni of the Indian Institute of Technology (BHU) Varanasi. My memory of study period at IIT (BHU) can never be complete without mentioning my fellow research scholars. Special thanks to Dr. Nagendra Pratap Singh, Dr. Arvind Tiwari, Dr. Vibhav Prakash Singh, Dr. Kuntesh K Jani, Dr. Gargi Srivastava, Dr. Roshan Singh, Dr. Tribikram Pradhan, Mr. Sushant Pandey, Mr. Ashwini Singh, Mr. Nigmendra Pratap Yadav, Mr. Ankit Jaiswal, Mr. Santosh Kumar Tripathy, Ms. Pratishtha Verma, Mrs. Divya Singh for their great help and cooperation.

I extend special thanks to the non-teaching staff in the Department, particularly, Mr. Manoj Singh, Mr. Ravi Kumar Bharati, Mr. Prakhar Kumar, Mr. Ritesh Singh and Mr. Shubham Pandey for their consistent support. My parents, Late. Neelam Sharma and Shree Kuwar Prasad Sharma, who gave me the power and brain to work out on this research and their blessings at every level made me to see this success. Words are insufficient to express my profound sense of gratitude to my loving wife Mrs. Khushaboo Sharma, as well as my loving parents-in-law Smt. Champak Lata Rai and Shree Ram Anant Rai whose encouragement and blessings gave me physical and moral strength throughout my career as well in the present research. I extend my thanks to my brothers Mr. Devesh Kumar Sharma and Mr Ashutosh Rai, and my sisters Mrs. Vibha Singh and Mrs. Ankita Singh, my great daughter Shaivya Sharma who is my parts of inspiration. Finally, I would like to wind up by paying my heartfelt thanks and prayers to the Almighty for his unbound love and grace.

Ashish Kumar Sharma

TABLE OF CONTENT

CERTIFICATE.....	iii
DECLARATION BY THE CANDIDATE	v
COPYRIGHT TRANSFER CERTIFICATE.....	vii
ACKNOWLEDGEMENTS.....	ix
Table of Content	xi
LIST OF FIGURES	xv
LIST OF TABLES.....	xvii
LIST OF ABBREVIATIONS.....	xix
LIST OF SYMBOLS	xxi
PREFACE.....	xxiii
Chapter 1: Introduction.....	1
1.1 Background	1
1.2 Motivation	3
1.3 Problem Statement	4
1.4 Research Objectives	5
1.5 Research Contributions	6
1.6 Outline of the Thesis	7
Chapter 2: Theoretical Background and Literature Review.....	9
2.1 Protein Secondary Structure Prediction	9
2.2 Deep Neural Network	12
2.2.1 Recurrent Neural Network	12
2.2.2 Gated Recurrent Unit.....	14
2.2.3 Long Short Term Memory.....	15
2.3 Research gaps and Findings:.....	17
2.4 Benchmark Datasets:.....	18
2.5 Performance Metrics:	18
2.6 Conclusions.....	20
Chapter 3: Protein Secondary Structure Prediction using Sequence Embedding and Bi-LSTM.....	21
3.1 Introduction	21

3.2 Protein Sequence Representation	25
3.3 Dataset Description.....	26
3.4 Model Architecture.....	27
3.4.1 Feature Extraction.....	27
3.4.2 Deep Network: Bidirectional Long Short Term Memory.....	27
3.4.3 Fully Connected Layer.....	29
3.5 Experimental Analysis.....	30
3.6 Results and Discussion	31
3.7 Conclusion.....	35
Chapter 4: Protein Secondary Structure Prediction using Character n-gram Embedding and Bi-LSTM	36
4.1 Introduction	36
4.2 Dataset Description.....	41
4.3 Problem Statement:.....	41
4.3.1 Input features	41
4.3.2 Outputs:.....	42
4.4 Methodology.....	43
4.4.1 Feature Extraction.....	43
4.4.2 Deep Network: Bidirectional LSTM	43
4.4.3 Fully Connected Layer.....	45
4.5 Experimental Details	45
4.5.1 Experimental Setup.....	45
4.5.2 Hyper-parameters.....	46
4.5.3 Training and Test Strategy.....	46
4.6 Results and Discussion	47
4.7 Conclusion.....	49
Chapter 5: Variable Length Character N-Gram Embedding Of Protein Sequences for Secondary Structure Prediction.....	50
5.1 Introduction	50
5.2 Protein Sequence Representation Method.....	53
5.3 Dataset Description.....	54
5.4 Methodology.....	55

5.4.1 Model Details	55
5.4.2 Character N-Gram	56
5.4.3 Deep Network: Bidirectional LSTM	57
5.4.4 Fully Connected Layer	57
5.5 Experimental Details	58
5.6 Results and Discussion.....	59
5.7 Conclusion.....	64
Chapter 6: Sequence-To-Sequence Modeling for Protein Secondary Structure Prediction	65
6.1 Introduction	65
6.2 Datasets	68
6.3 Proposed Model	72
6.4 Results and Discussion.....	74
6.5 Conclusion.....	76
Chapter 7: Conclusion and Future Work.....	79
7.1 Conclusions	79
7.2 Future Work	82
List of Papers Published /Communicated.....	95

LIST OF FIGURES

Figure 2.1: Recurrent Neural Network	13
Figure 2.2: Gated Recurrent Unit	14
Figure 2.3: Long Short Term Memory Network	16
Figure 3.1: Model Architecture of Protein Secondary Structure Prediction.....	29
Figure 3.2: Plot of Training Accuracy and testing Accuracy over Epochs for CullPDB32	
Figure 3.3: Plot of Training Accuracy and testing Accuracy over Epochs for SS.txt....	32
Figure 3.4: plot for Training and Testing Loss over Epochs for SS.txt.....	33
Figure 3.5: plot for Training and Testing Loss over Epochs for CullPDB.....	33
Figure 4.1: Model Architecture of Protein Secondary Structure Prediction.....	44
Figure 4.2: Plot for Training and Testing Accuracy over Epochs for ss.txt	47
Figure 4.3: Plot for Training and Testing Loss over Epochs for ss.txt.....	47
Figure 5.1: Proposed Model for Protein Secondary Structure Prediction	56
Figure 5.2: Training and Testing Loss over Epochs for ss.txt.....	60
Figure 5.3: Training and Testing Accuracy over Epochs for ss.txt	60
Figure: 6.1 Sequence to Sequence using LSTM [8]	69
Figure 6.2: Sequence to sequence model for protein secondary structure prediction	72

LIST OF TABLES

Table 3.1: DSSP 8-class to 3-class Conversion	26
Table 3.2: Configuration of Training Hyper-parameters	31
Table 3.3: Performance of Bi-RNN, Bi-GRU and Bi-LSTM on ss.txt dataset	32
Table 3.4: Comparison of the performance of various	34
Table 3.5: Performance Accuracy comparison with other Methods on Dataset RS126.	35
Table 4.1: DSSP 8-class to 3-class conversion.....	42
Table 4.2: Configuration of Training Hyper-parameters	46
Table 4.3: Performance Q3 Accuracy comparison with other methods on CASp9 dataset	48
Table 4.4: Performance Q3 Accuracy comparison with other methods on RS126 dataset.	49
Table 5.1: Hyper-parameters Value.....	58
Table 5.2: Performance Comparison of Bi-RNN, Bi-GRU and Bi-LSTM	59
Table 5.3: Q3 accuracy for different value of n-gram	61
Table 5.4: Q3 accuracy for different combination of n-gram.....	62
Table 5.5: Q3 value comparison with other Methods on dataset CASP9.....	63
Table 5.6: Q3 value comparison with other Methods on dataset RS126.....	63
Table 6.1: Amino acid one-hot Encoding Representation.....	69
Table 6.2: Comparison of the performance of various single-sequence based prediction on Cullpdb.....	75
Table 6.3: Comparison of various single-sequence based prediction on data1199	76

LIST OF ABBREVIATIONS

AI	Artificial Intelligence
ASA	Accessible Surface Area
Bi-GRU	Bidirectional Gated Recurrent Unit
Bi-LSTM	Bidirectional Long Short Term Memory
Bi-RNN	Bidirectional Recurrent Neural Network
CNN	Convolutional Neural Network
CPU	Central Processing Unit
DL	Deep Learning
DNN	Deep Neural Network
DSSP	Define Secondary Structure of Protein
GOR	Garnier Osguthorpe-Robson
GPU	Graphics Processing Unit
GRU	Gated Recurrent Unit
HMM	Hidden Markov Model
LSTM	Long Short Term Memory
ML	Machine Learning
MSA	Multiple Sequence Alignment
NLP	Natural Language Processing
NMR	Nuclear Magnetic Resonance
NN	Neural Network
PSSM	Position Specific Scoring Matrix
TPR	True Positive Rate

TNR	True Negative Rate
FPR	False Positive Rate
FNR	False Negative Rate
PSSP	Protein Secondary Structure Prediction
RNN	Recurrent Neural Network
SVM	Support Vector Machine
RF	Random Forest
PseAAC	Pseudo Acid Composition

LIST OF SYMBOLS

B	Beta Bridge
b	biases
C	Coil Regions
E	Beta Strand
F	Activation Function
f_t	Forget Gate
G	310-helix
H	Alpha Helix
h_t	State Vector
I	pi-helix
i_t	Input Gate
K	Secondary Structure Classes
L	Length of Amino Acid Sequence
o_t	Output Gate
Q3	Three State Per Residue Accuracy
Q8	eight State Per Residue Accuracy
S	loop
T	Beta Turn
tanh	Tangent Function
V	Size of n-gram words vocabulary
W	Weight Matrix
x_t	Input Vector

Z	Probability Vector
\odot	Hadamard Product
Σ	Sigmoid Function

PREFACE

Protein is an essential molecule for defining the structure and function of body cells. However, protein sequences have enough information for determining their structure and function. There is a large gap between the number of protein sequences and the number of protein structures. Protein secondary structure prediction from its amino acid sequence is an important problem in molecular biology. Protein structural information is critical for understanding their function, but many protein sequences do not have their structural information. The large gap between protein sequence and structure motivates computational methods and, in particular, deep learning methods predicting structural information from the protein primary sequences.

In this thesis, we focused on predicting protein secondary structure from the primary sequences. Protein primary sequences represented using subsequence-based representations of amino acid sequences. The dense embedding vectors of primary sequences are trained over a large dataset of primary sequences. These dense embedding vectors combined with a bidirectional long short term memory network to capture the contextual information of amino acid residue for predicting secondary structure.

Firstly the thesis work focuses on primary sequence representation for deep learning methods. The primary sequences are segmented into twenty amino acid characters.

The dense embedding vector of amino acid residues combined with deep learning techniques for protein secondary structure prediction. Further, the primary sequences are a bag of character n-gram words. The different sizes of n-gram words computed and the dense vector of character n-gram words used to capture the local contextual

information for secondary structure prediction. The variable-length character n-gram words dense embedding vector used for secondary structure prediction to overcome the drawback of fixed size character n-gram word. The recurrent neural network used for secondary structure prediction, the bidirectional long short-term memory network, better capture the long contextual information.

Further protein secondary structure prediction mapped as a sequence to sequence modeling. The amino acid residues are labeled with secondary structure elements. The state vectors for encoder and decoder are computed using long short term memory network.