



NAMED ENTITY BOUNDARY DETECTION FOR SINHALA

Yapa Hetti Pathirannahalage Prasan Priyadarshana

208073P

Thesis/Dissertation submitted in partial fulfillment of the requirements for the degree
Master of Science (Major Component Research)

Department of Information Technology
Faculty of Information Technology

University of Moratuwa
Sri Lanka

March 2022

Declaration

I declare that this is my own research thesis, and this thesis does not incorporate without acknowledgement any material previously published submitted for a Degree or Diploma in any other university or institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text.

Signature: *UOM Verified Signature*

Date: 18/03/2022

I have read the final thesis and it is in accordance with the approved university final thesis outline.

Signature of the supervisor: *UOM Verified Signature*

Date: 18/03/2022

Acknowledgement

I would take this opportunity to express my sincere gratitude to everyone who guided me directly and indirectly to achieve this milestone. First and foremost, I am extremely grateful to my supervisor, Dr. Lochandaka Ranathunga for the continuous support, invaluable advice, kind guidance and patience during this whole journey. His vast experience and extensive knowledge have encouraged me throughout the time of my academic research and life. I also thank the research assistants who have been employed under the Social Media Analytics Research Group and the Accelerating Higher Education Expansion and Development (AHEAD) project, University of Moratuwa who helped me in annotating the dataset. I would also like to thank Dr. Surangika Ranathunga and all the research assistants who have been attached to the National Languages Processing Centre (NLPC), University of Moratuwa for providing me various corpus related to named entities and shaping the overall idea to identify the novelty. Finally, I would like to express my gratitude to my beloved parents, my wife and my friends for their motivation, inspiration, and continuous support throughout the difficult times.

Abstract

Named entity recognition (NER) can be introduced as a sequential categorizing task which contains a potential gravity in novel research arena. NER can be mentioned as the foundation for accomplishing most common natural language processing (NLP) tasks such as information extraction, information retrieval, semantic role labelling etc. Even though plenty of attempts have been employed on NE type detection, still there are plenty of avenues to be discovered under the NE boundary detection. Analyzing Sinhala related contents which have been published in social media can also be considered as one of the rising factors due to several human involvements in the recent past. The ultimate goal which is to obtain a constructive deep neural framework for determining named entity boundary detection has been achieved in a comprehensive manner and the model has been tested using Sinhala related statements which have been extracted through social media. Several objectives have been determined to accomplish this task considering the existing baselines. Several novelties have been identified to show off the uniqueness of this approach. Specifically, the novel concept “*Boundary Bubbles*” has been used to identify the specific entity mentions considering each head word for the identified named entities. Various experiments have been conducted based on multiple evaluation criteria and the named entity boundary detection model performs well with an average of 71% in Precision, 67% in Recall and 63% in F1 over the existing benchmarks. Hence this novel framework can be accepted as a vital solution for determining named entity boundary detection under forecasting various computational activities in social media.

Keywords:

Deep neural network, named entities, named entity boundary, named entity recognition, named entity type

TABLE OF CONTENTS

Declaration of the Candidate and Supervisor	ii
Acknowledgement	iii
Abstract	iv
Table of Contents	v
List of Figures	vii
List of Tables	vii
List of Abbreviations	viii
1. Introduction	1
1.1 Problem Domain	1
1.2 Background	3
1.3 Motivation	3
1.4 Problem Statement	4
1.5 Research Aim & Objectives	4
2. Literature Review	8
2.1 Literature review on NE boundary detection	9
2.2 Analysis of Sinhala social media statements	17
3. Research Methodology	23
3.1 Methodological Approach for Sinhala Social Media Fundamentals	23
3.2 Methodological Approach for Deriving NE Boundary Detection	29
4. Implementation	40
4.1 Technology Stack	40
4.2 Implementation on Sinhala NE Identification	41
4.3 Implementation on Deep Transfer Learning for Word Representation	42
4.4 Implementation on NE Mention Head Word Detection	44
4.5 Implementation on the Mention Nuggets Identifier and The Head Detection Loss Function	47
4.6 Implementation on NE Boundary Detection Regional Classification	50
4.7 Implementation on NE Linking	52

5. Experiments	53
5.1 Experimental Settings	53
5.1.1 Dataset	53
5.1.2 Tools and Libraries	54
5.1.3 Hardware Requirements	54
5.2 Experimental Baselines	55
5.3 Experimental Results	56
5.3.1 Testing on Sinhala NE Identification	56
5.3.2 Testing on word embedding through deep transfer learning	58
5.3.3 Testing on NE head word detection	59
5.3.4 Testing on NE boundary detection	61
5.3.5 Testing on NE linking	61
6. Evaluation & Discussion	63
6.1 Evaluation on Sinhala NE identification	63
6.2 Evaluation on word embedding through deep transfer learning	66
6.3 Evaluation on NE head word detection	68
6.4 Evaluation on NE boundary detection	70
6.5 Evaluation on NE linking	71
7. Conclusion	73
References	76

LIST OF FIGURES

Figure	Description	Page
Figure 1.1	Hierarchy of NER	2
Figure 3.1	Pseudocode for POS tagging	28
Figure 3.2	Proposed architecture for named entity boundary detection	29
Figure 3.3	Deep Transfer Learning process	30
Figure 3.4	Main system architecture of NE type detection	32
Figure 3.5	Abstract representation of boundary bubbles	33
Figure 3.6	Stack LSTM architecture of NE boundary detection	37
Figure 6.1	Performance comparison on deep transfer learning	66
Figure 6.2	GloVe & XLM-R performance comparison on transfer learning	67

LIST OF TABLES

Table	Description	Page
Table 3.1	Bag of Extracted Stop Words	25
Table 3.2	Stems and Variations	26
Table 5.1	Results on Sinhala NE Detection	57
Table 5.2	Hyper Parameter Based Results on Sinhala NE Detection	57
Table 5.3	Average Results on Word Embeddings	59
Table 5.4	Average Results on NE Head Word Detection	59
Table 5.5	Parameters on Mention Nuggets Identifier	60
Table 5.6	Value Adjustments on Mention Nuggets Identifier	60
Table 5.7	Average Results on NE Boundary Detection	61
Table 5.8	Results on NE Linking	62
Table 5.9	Average Results on NE Linking	62
Table 6.1	Kappa Values for Individual NE Categories	64
Table 6.2	Overall Kappa Values for NE Categories	64
Table 6.3	Evaluation on NE Type Classification	65
Table 6.4	Evaluation on NE Head Word Detection	68
Table 6.5	Summary Evaluation on NE Head Word Detection	70
Table 6.6	Evaluation on NE Boundary Detection	70
Table 6.7	Evaluation on NE Linking	71

LIST OF ABBREVIATIONS

Abbreviation	Description
NER	Named Entity Recognition
NE	Named Entities
NLP	Natural Language Processing
ORG, PER, LOC	Identified Some of the Named Entities
CNN	Convolutional Neural Networks
LSTM	Long Short-Term Memory
Bi-LSTM	Bidirectional Long Short-Term Memory
GPU	Graphics Processing Unit
GRU	Gated Recurrent Units
Bi-GRU	Bidirectional Gated Recurrent Units
BERT	Bidirectional Encoder Representations from Transformers
MRC	Machine Reading Comprehensive
BOW	Bag of Words
RNN	Recurrent Neural Network
MLP	Multi-Layer Perceptron
CRF	Conditional Random Fields
NED	Named Entity Disambiguation
IE	Information Extraction
IR	Information Retrieval