

**Pre-training and Fine-tuning Multilingual
Sequence-To-Sequence Models for Domain-Specific
Low-Resource Neural Machine Translation**

Sarubi Thillainathan

208037K

Thesis/Dissertation submitted in partial fulfillment of the requirements for the degree
Master of Science in Computer Science and Engineering

Department of Computer Science & Engineering

University of Moratuwa

Sri Lanka

April 2022

DECLARATION

I, Sarubi Thillainathan, declare that this is my own work and this dissertation does not incorporate without acknowledgement any material previously 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.

Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my dissertation, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works (such as articles or books).

Signature:

Date:

The above candidate has carried out research for the Masters thesis/Dissertation under my supervision.

Name of Supervisor: Dr. Surangika Ranathunga

Signature of the Supervisor:

Date:

Name of Supervisor: Prof. Sanath Jayasena

Signature of the Supervisor:

Date:

ABSTRACT

Limited parallel data is a major bottleneck for morphologically rich Low-Resource Languages (LRLs), resulting in Neural Machine Translation (NMT) systems of poor quality. Language representation learning in a self-supervised sequence-to-sequence fashion has become a new paradigm that utilizes the largely available monolingual data and alleviates the parallel data scarcity issue in NMT. The language pairs supported by the Self-supervised Multilingual Sequence-to-sequence Pre-trained (SMSP) model can be fine-tuned using this pre-trained model with a small amount of parallel data.

This study shows the viability of fine-tuning such SMSP models for an extremely low-resource domain-specific NMT setting. We choose one such pre-trained model: mBART. We are the first to implement and demonstrate the viability of non-English centric complete fine-tuning on SMSP models. To demonstrate, we select Sinhala, Tamil and English languages in extremely low-resource settings in the domain of official government documents.

This research explores the ways to extend SMSP models to adapt to new domains and improve the fine-tuning process of SMSP models to obtain a high-quality translation in an extremely low-resource setting. We propose two novel approaches: (1) Continual Pre-training of the SMSP model in a self-supervised manner with domain-specific monolingual data to incorporate new domains and (2) multistage fine-tuning of the SMSP model with in- and out-domain parallel data.

Our experiments with Sinhala (Si), Tamil (Ta) and English (En) show that directly fine-tuning (single-step) the SMSP model mBART for LRLs significantly outperforms state-of-the-art Transformer based NMT models in all language pairs in all six bilingual directions. We gain a +7.17 BLEU score on Si→En translation and a +6.74 BLEU score for the Ta→En direction. Most importantly, for non-English centric Si-Ta fine-tuning, we surpassed the state-of-the-art Transformer based NMT model by gaining a +4.11 BLEU score on Ta→Si and a +2.78 BLEU score on Si→Ta.

Moreover, our proposed approaches improved performance strongly by around a +1 BLEU score compared to the strong single-step direct mBART fine-tuning for all six directions. At last, we propose a multi-model ensemble that improved the performance in all the cases where we obtained the overall best model with a +2 BLEU score improvement.

Keywords: Neural Machine Translation, Pre-trained Language Models, Pre-training, Fine-tuning, Low-Resource languages, mBART

DEDICATION

With deepest gratitude, I dedicate this research to my Grandpa, *Late M. Saravanamuttu*.
And, of course, to my forever loving family for supporting my dreams, no matter what!

ACKNOWLEDGEMENTS

First and foremost, I would like to express my sincere gratitude to my supervisors Dr Surangika Ranathunga and Prof. Sanath Jayasena, for the continuous support and guidance given for the success of this research. Without your insights and tremendous mentorship, I could not have achieved this level. Thank you for always stepping in to help whenever I needed regardless of your busy schedules.

I wish to convey my sincere appreciation to Prof. Gihan Dias and National Language Processing Center (NLPC) members for their valuable insights and support given for this research. I would also like to thank the entire Department of Computer Science and Engineering staff, both academic and non-academic, for their help and for providing me with the resources necessary to conduct my research. This research was supported by the University of Moratuwa AHEAD project Research Grant.

Lastly, I want to thank my family and friends who supported me though this journey.

Thank you!

TABLE OF CONTENTS

Declaration of the Candidate & Supervisor	i
Abstract	ii
Dedication	iii
Acknowledgement	iv
Table of Contents	v
List of Figures	viii
List of Tables	ix
List of Abbreviations	1
1 Introduction	2
1.1 Background	2
1.2 Research Problem	3
1.3 Research Scope and Objectives	4
1.4 Contributions	4
1.5 Publications	5
2 Literature Survey	6
2.1 Overview	6
2.2 Neural Machine Translation (NMT)	6
2.3 NMT for Sinhala, Tamil, and English Languages	7
2.4 Multilingual Neural Machine Translation (MNMT)	8
2.4.1 Universal Encoder and Decoder Architecture for MNMT	10
2.4.2 Strategies to improve MNMT	10
2.5 Transfer Learning (TL) in NMT	11
2.5.1 Fine-tuning techniques	12
2.5.2 Transfer Learning Protocols	13
2.6 Pre-trained Models for NMT	14
2.7 Self-supervised Multilingual Sequence-to-sequence Pre-training	15
2.7.1 BART	16
2.7.2 mBART	16

2.8	Fine-tuning Multilingual Self-Supervised Pre-trained Models for Low-resource NMT	17
2.9	Continual learning on Self-Supervised Pre-trained Models (Extending the pre-trained models)	19
2.10	Summary	20
3	Methodology	22
3.0.1	Overview	22
3.0.2	Bilingual Fine-tuning using Self-Supervised Multilingual Sequence-to-sequence Pre-trained Models	23
3.0.3	Multilingual Fine-Tuning using Self-Supervised Multilingual Sequence-to-sequence Pre-trained Models	24
3.1	Continual Pre-training for Domain Adaptation	24
3.2	Multistage Fine-tuning	26
3.2.1	Two-stage FT	28
3.2.2	Multistage Fine-tuning Combine with Continual Pre-Training	28
3.2.3	Ensemble of Fine-tuned Self-Supervised Multilingual Sequence-to-sequence Pre-trained Models	29
4	Implementation	31
4.1	Experimental setup	31
4.1.1	Architecture	31
4.1.2	Dataset	31
4.1.3	Preprocessing	32
4.2	Addressing the Zero Width Joiner (ZWJ) issue	33
4.3	Baselines	34
4.4	Fine-tuning SMSP Model	34
4.5	Continual Pre-training for Domain Adaptation	35
4.6	Multistage Fine-tuning	36
4.7	Evaluation setup	36
5	Results and Discussion	37
5.1	Bilingual Fine-tuning using Multilingual Denoising Pre-trained Models	37
5.2	Multilingual Fine-tuning using Multilingual Denoising Pre-trained Models	39

5.3	Continual Pre-training for Domain Adaptation	39
5.4	Multistage Fine-tuning	40
5.4.1	Two-stage FT	40
5.4.2	Multistage fine-tuning Combined with Continual Pre-Training	41
5.5	Ensemble of Fine-tuned Self-supervised Multilingual Sequence-to-sequence Pre-trained Models	42
5.6	Manual Analysis of the Translated Output	43
6	Conclusion and Future work	44
	References	45
A	Appendix	56
A.0.1	Addressing the Zero Width Joiner (ZWJ) issue	56
A.0.2	Output Translated Sentences	56

LIST OF FIGURES

Figure 2.1	Overview of MNMT Categories [1]	8
Figure 2.2	Overview of MNMT Architectures [1]	9
Figure 2.3	Overview of Transfer Learning	12
Figure 2.4	Overview of Multilingual Denoising Autoencoder Pre-training (left) and fine-tuning on NMT (right) [2]. A special token "language id" is added to both the encoder and decoder.	17
Figure 2.5	Overview of Pre-training and fine-tuning.	18
Figure 3.1	Overview of Methodology	22
Figure 3.2	Overview of Continual Pre-training for Domain Adaptation.	25
Figure 3.3	Overview of Multistage Fine-tuning.	27
Figure 3.4	Different ways of Multistage Fine-tuning.	27
Figure A.1	Output Translated Sentences	57

LIST OF TABLES

Table 2.1	Sample input and its transformations output after applying different noising functions [3]	16
Table 4.1	Statistics of the parallel dataset of official government documents	32
Table 4.2	Statistics of the out-domain parallel corpus	32
Table 4.3	Monolingual Data	33
Table 5.1	Comparison between full precision training and mixed precision Fine-Tuning. Results are reported in BLEU score.	37
Table 5.2	Comparison with SMT, LSTM, Transformer Architectures against our Bilingual Fine-tuning models for Sinhala (Si), Tamil (Ta) and English (En) - Results are reported in BLEU score.	38
Table 5.3	Results of Bilingual Fine-tuning models and Multilingual Sinhala Centric Fine-tuning models - Results are reported in BLEU score.	39
Table 5.4	Fine-tuning Results from Continual Pre-trained models against our strong baseline Bilingual Fine-tuned models for Si↔En pairs - Results are reported in BLEU score.	40
Table 5.5	Fine-tuning Results of Bilingual and Trilingual Continual Pre-training on in-domain monolingual data for all the six directions - Results are reported in BLEU score.	40
Table 5.6	Fine-tuning Results from Continual Pre-trained models against the our strong baseline Bilingual Fine-tuned models for Si↔En pairs - Results are reported in BLEU score.	41
Table 5.7	Multistage fine-tuning against the our strong baseline Bilingual Fine-tuned models - Results are reported in BLEU score.	42
Table 5.8	Top 4 improved models from baseline B-FT	42
Table 5.9	Ensemble Results for all the six directions. Results are reported in BLEU score.	43

LIST OF ABBREVIATIONS

NLP	Natural Language Processing
MT	Machine Translation
SMT	Statistical Machine Translation
NMT	Neural Machine Translation
MNMT	Multilingual Neural Machine Translation
RNN	Recurrent Neural Networks
LSTM	Long Short Term Memory
GRU	Gated Recurrent Unit
DAE	Denoising Autoencoder
MLE	Maximum Likelihood Estimate
TL	Transfer Learning
SMSP	Self-supervised Multilingual Sequence-to-sequence Pre-trained
FT	Fine-Tuning
LRL	Low-Resource Language
LM	Language Model
M-FT	Multilingual Fine-Tuning
B-FT	Bilingual Fine-Tuning
CPT	Continual Pre-Training