

CERTIFICATE

It is certified that the work contained in the thesis titled "*Low Resource Similar Language Machine Translation in Common Phonetic Space with Multilingual, Adversarial and Reinforced Deep Learning*" by *Amit Kumar* 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 requirements of Comprehensive Examination, Candidacy, and SOTA for the award of Ph.D. Degree.



Supervisor

Dr. Anil Kumar Singh

Associate Professor,

Dept. of Computer Science and Engineering,

Indian Institute of Technology (BHU),

Varanasi, India 221005.

पर्यवेक्षक/Supervisor
संगणक विज्ञान एवं अभियंत्रण विभाग
Department of Computer Sc. & Engg
भारतीय प्रौद्योगिकी संस्थान
Indian Institute of Technology
(काशी हिन्दू विश्वविद्यालय)
Banaras Hindu University
वाराणसी/Varanasi-221005



Co-Supervisor

Dr. Ajay Pratap

Assistant Professor,

Dept. of Computer Science and Engineering

Indian Institute of Technology (BHU),

Varanasi, India 221005.

सहायक आचार्य/Assistant Professor
संगणक विज्ञान एवं अभियंत्रण विभाग/Department of Computer Sc. & Engg
भारतीय प्रौद्योगिकी संस्थान/Indian Institute of Technology
(काशी हिन्दू यूनिवर्सिटी)/(Banaras Hindu University,
वाराणसी /Varanasi-221005

DECLARATION BY THE CANDIDATE

I, **Amit Kumar**, certify that the work embodied in this Ph.D. thesis is my own bonafide work carried out by me under the supervision of **Dr. Anil Kumar Singh** and co-supervision of **Dr. Ajay Pratap** from **July 2017** to **December 2022** at **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 copied any other's work, paragraphs, text, data, results, *etc.* reported in journals, books, magazines, reports, dissertations, theses, *etc.*, or available at websites and have not included them in this thesis and have not cited as my own work.

Date: 06/12/2022

Place: Varanasi


(Amit Kumar)

CERTIFICATE BY THE SUPERVISOR

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.



(Supervisor)

Dr. Anil Kumar Singh

Associate Professor,

पर्यवेक्षक (Supervisor)
समगणक विज्ञान एवं अभियांत्रिकी विभाग
Department of Computer Sc. & Engg
भारतीय प्रौद्योगिकी संस्थान
Indian Institute of Technology
(काशी हिन्दू विश्वविद्यालय)
(Banaras Hindu University)
वाराणसी/Varanasi-221005



(Co-Supervisor)

Dr. Ajay Pratap

Assistant Professor,

सहायक आचार्य/Assistant Professor
समगणक विज्ञान एवं अभियांत्रिकी विभाग
Department of Computer Sc. & Engg
भारतीय प्रौद्योगिकी संस्थान/Indian Institute of Technology
(बनारस हिन्दू यूनिवर्सिटी)/(Banaras Hindu University)
वाराणसी/Varanasi-221005


Signature of Head of Department

Professor & Head

समगणक विज्ञान एवं अभियांत्रिकी विभाग
Department of Computer Sc. & Engg
भारतीय प्रौद्योगिकी संस्थान
Indian Institute of Technology
(बनारस हिन्दू यूनिवर्सिटी)
(Banaras Hindu University)
वाराणसी-221005

COPYRIGHT TRANSFER CERTIFICATE

Title of the Thesis: Low Resource Similar Language Machine Translation in Common Phonetic Space with Multilingual, Adversarial and Reinforced Deep Learning

Name of the Student: Amit Kumar

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: 06/12/2022

Place: Varanasi



(Amit Kumar)

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.

Dedicated to my parents,

Shri. Panna Lal

and

Smt. Bhagirathi Devi

ACKNOWLEDGEMENT

First and foremost, I would like to thank my supervisor, Dr. Anil Kumar Singh, and my co-supervisor, Dr. Ajay Pratap for their invaluable support and assistance. I feel immense pleasure in expressing my profound sense of gratitude and sincere regard for their constant feedback and expertise during all these years. I am eternally grateful to have the opportunity to work on my thesis under their supervision.

My cordial thanks to all the members of the Department of Computer Science and Engineering for creating an excellent working atmosphere. I would also like to thank the other members of my Doctoral committee, Dr. Sukhda, Department of Humanities Studies, and Dr. Sukomal Pal, Department of Computer Science and Engineering, for their help and support throughout the tenure of my studies. Special thanks to Dr. Sriparna Saha, Department of Computer Science and Engineering, IIT Patna, for her valuable suggestions. I would also like to convey my sincere gratitude to Prof. Sanjay Kumar Singh, Head of the CSE Department, and all the RPEC and DPGC members for their suggestions and endorsement of this work.

I am grateful to my colleagues and friends for the long discussions and their brilliant insights that have helped me to overcome the challenges I have faced in the development of this work. Finally, I express my heartfelt gratitude to my parents Smt. Bhagirathi Devi and Shri. Panna Lal for their constant support, love, encouragement, and sacrifices. Their affectionate love and care cannot be expressed in words.

With limitless humility, I would like to praise and thank the “ **Satguru Kabir Sahab Ji**”. The almighty, the Merciful compassionate who bestowed me with all the favourable circumstances to achieve the desired goal of life through this crucial juncture.

(Amit Kumar)

List of Figures

1.1	Thesis organization.	10
2.1	Reinforcement Learning.	17
3.1	Proposed architecture.	36
4.1	Zero Resource Problem in MT	54
4.2	TSL for zero-resource language pair.	57
4.3	MSL for zero-resource language pair.	59
4.4	TLSPG approach.	60
5.1	Examples of sentences in different languages.	74
5.2	Weight initialization by language model.	77
5.3	Architecture of RSSW.	79
5.4	Changing the amount of the selected pseudo in-domain data on the HI \leftrightarrow NE task.	88
5.5	BLEU score versus Gaussian iteration for HI \leftrightarrow NE task.	88
6.1	Some examples of IPA.	92
6.2	Joint embedding.	98
6.3	Optimized generator model in GAN-NMT.	99
6.4	Overall architecture of GAN NMT.	102
6.5	Flow chart of proposed training procedure.	103
6.6	chrF2 scores.	108
6.7	TER scores.	108
6.8	Translation examples by different systems.	110

List of Tables

2.1	Comparative overview of the closely related existing models for zero-shot problem	19
2.2	Comparison of existing works for morphological richness issue. ✓ and ✗ represent presence and absence of particular feature, respectively.	21
2.3	Comparison of related work for domain shift problem	23
2.4	Comparison of existing works for context missing and rare-word problem	26
3.1	Some details about the languages used in our experiments	34
3.2	Corpus Statistics showing the number of training, validation, and test sentences for each domain	38
3.3	Experiment results (BLEU, chrF2, and TER scores).	39
3.4	BLEU score-based comparison of SMT, SMT + WX and the proposed approaches.	41
3.5	Similarity between languages using SSNGLMScore	41
3.6	char-BLEU score on the training data	43
3.7	TER and chrF2 scores on the training data	44
3.8	Character-based entropy of languages with or without applying WX-notation	46
3.9	Character-level Entropy computed on Monolingual Vocabulary	46
3.10	Character-level Redundancy Reduction on Monolingual Corpus	46
3.11	Character-level Entropy and Redundancy Changes for Parallel Corpus	47
3.12	Cross-lingual distance between languages after applying character-level language model using perplexity-based score (Unnormalized on language directions)	48
3.13	Cross-lingual distance between languages after applying character-level language model using perplexity-based score	48
3.14	Experiments on distant language pairs.	49
3.15	Applying on zero-shot language pairs.	50

3.16	Experiments on back-translation.	50
4.1	Relatedness features between languages	55
4.2	Description of corpus statistics	63
4.3	Experimental setup used to train the TSL model	64
4.4	Experimental results for different language pairs	65
4.5	SSNGLMScore	68
4.6	Entropy	70
4.7	Experiments on using LSTM instead of Transformer	70
5.1	Corpus statistics	84
5.2	Results on HI \leftrightarrow NE translation task	86
5.3	Results on HI \leftrightarrow MR translation task	86
5.4	Comparison with state-of-the-art approach on HI \leftrightarrow NE translation task	89
5.5	Results of LSTM as translation model on HI \leftrightarrow NE translation task	90
6.1	Results on different methods (“ \rightarrow ”: X \rightarrow HI and “ \leftarrow ”: HI \rightarrow X)	105
6.2	Results on different methods under multilingual settings (“ \rightarrow ”: X \rightarrow HI and “ \leftarrow ”: HI \rightarrow X)	105
6.3	Unsupervised results on X \rightarrow HI and HI \rightarrow X	107
7.1	Overview of the proposed approaches	111
7.2	Summarization of Chapter-3	112
7.3	Summarization of Chapter-4	112
7.4	Summarization of Chapter-5	112
7.5	Summarization of Chapter-6	113

List of Symbols

Symbol	Description
X	Source sequence
Y	Target sequence
x_i	Source word embedding at i^{th} position in X
y_j	Word at position j in target sequence
q_t	Hidden state at time t
α_{ti}	Attention weight between source word x_i and target word y_t
$output_length$	Length of predicted sentences
$reference_length$	Length of reference sentences
$attn$	Attention in Transformer
Q	Query
K	Key
V	Value
d_k	Dimension of key
\mathcal{L}_{CE}	Cross-entropy loss
\mathcal{L}_{LM}	Cross-entropy loss of language model
PP	Perplexity
w_i	i^{th} word embedding
W_i	initial weight of each sentence pair $S_i (X_i, Y_i)$
c	State
a	Action
θ	Models parameter in policy network
$\pi_{\theta}(c, a)$	Policy
\mathcal{N}	Gaussian distribution
μ	Mean
σ	Standard deviation
R	Reward

Symbol	Description
sw_i	Subword embeddings at i^{th} position
so_i	Sub-phonetics embeddings at i^{th} position

Abbreviations

Abbreviation	Description
NLP	Natural Language Processing
MT	Machine Translation
NMT	Neural Machine Translation
SMT	Statistical Machine Translation
CNN	Convolutional Neural Network
LSTM	Long-Short Term Memory
GAN	Generative Adversarial Network
NN	Neural Network
RNN	Recurrent Neural Network
GRU	Gated Recurrent Unit
ReLU	Rectified Linear Unit
IPA	International Phonetic Alphabet
HRLs	High Resource Languages
LRLs	Low Resource Languages
ZRP	Zero-Resource Problem
ZST	Zero-Shot Translation
TER	Translation Edit Rate
TL	Transfer Learning
DA	Domain Adaptation