

**Combining Automatic Speech Recognition Models To  
Reduce Error Propagation in Low-Resource Transfer-  
Learning Speech-Command Recognition**

Jazeem Mohamed Isham  
(209333X)

Degree of Master of Science

Department of Computer Science and Engineering  
Faculty of Engineering

University of Moratuwa  
Sri Lanka

March 2022

# **Combining Automatic Speech Recognition Models To Reduce Error Propagation in Low-Resource Transfer-Learning Speech-Command Recognition**

Jazeem Mohamed Isham

(209333X)

Thesis/Dissertation submitted in partial fulfillment of the requirements for the degree

Master of Science specializing in Data Science

Department of Computer Science and Engineering

Faculty of Engineering

University of Moratuwa

Sri Lanka

March 2022

## DECLARATION

I declare that this is my work and this dissertation does not incorporate without acknowledgment 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 acknowledgment is made in the text.

Also, I hereby grant to the University of Moratuwa the non-exclusive right to reproduce and distribute my dissertation, in whole or in part in print, electronic, or another 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 Master's dissertation under my supervision. I confirm that the declaration made above by the student is true and correct

Name of the supervisor: Dr. Uthayasanker Thayasivam

Signature of the supervisor:

Date:

## ABSTRACT

There are several applications when comes to spoken language understanding such as topic modeling and intent detection. One of the primary underlying components used in spoken language understanding studies is automatic speech-recognition models. In recent years we have seen a major improvement in the automatic speech recognition system to recognize spoken utterances. But it is still a challenging task for low-resource languages as it requires hundreds of hours of audio input to train an automatic speech recognition model.

To overcome this issue recent studies have used transfer learning techniques. However, the errors produced by the automatic speech recognition models significantly affect the downstream natural language understanding models used for intent or topic identification. In this work, we have proposed a multi-automatic speech recognition set up to overcome this issue. We have shown that combining outputs from multiple automatic speech recognition models can significantly increase the accuracy of low-resource speech-command transfer-learning tasks than using the output from a single automatic speech recognition model.

We have come up with convolution neural network-based setups that can utilize outputs from pre-trained automatic speech recognition models such as DeepSpeech2 and Wav2Vec 2.0. The experiment result shows a 7% increase in accuracy over the current state-of-the-art low resource speech-command phoneme-based speech intent classification methodology.

## **ACKNOWLEDGMENTS**

First of all, I would like to thank my supervisor Dr. Uthayasanker Thayasivam for being my supervisor and guiding me throughout the project. His expertise in this related area helped me a lot in setting the right direction for this project and obtaining the needed datasets.

I would not have achieved this without the immense support from my family. I would like to thank my parents for supporting me throughout my studies from the start the up until today. A special thanks to my wife and family for coping with me throughout my research program and encouraging me to complete it.

I also want to thank the University of Moratuwa for giving me an opportunity to participate in the MSc program and for providing the necessary resources to complete this research. This would not have been easy without the support from my workplace as well. I would like to thank Sysco LABS Srilanka for allowing me to do this part-time MSc and research while working with them. Also, my colleagues have been very understanding and allowed me to spend time on my research even during busy office schedules.

# TABLE OF CONTENTS

DECLARATION	i
ABSTRACT	ii
ACKNOWLEDGMENTS	iii
TABLE OF CONTENTS	iv
LIST OF FIGURES	vi
LIST OF TABLES	vii
LIST OF ABBREVIATIONS AND ACRONYMS	viii
1 INTRODUCTION	1
1.1 Background	1
1.2 Research Problem	2
1.3 Research Objectives	3
1.4 Outline	3
2 LITERATURE REVIEW	4
2.1 Low-resource Transfer Learning	4
2.2 DeepSpeech2	5
2.3 Wav2Vec 2.0	6
2.4 Dual-Input CNN Model	7
2.5 Feature Concatenation	7
3 METHODOLOGY	9
3.1 Multi-ASR Combinations	9
3.2 Method 1 - Dual-input CNN Models	10
3.2.1 Input Layer (DeepSpeech2 and Wav2Vec 2.0)	10
3.3 Method 2 - Feature Concatenation	12
3.4 Commonly Used Techniques In Both Models	15

3.4.1	Convolutional Layer	15
3.4.2	Activation Function	15
3.4.3	Max-pooling and Dropout Layers	15
4	DATA SET	16
5	EXPERIMENT	18
5.1	Hyperparameter Tuning	18
5.2	K-Fold Cross-Validation	19
5.3	Experiment Setup	20
6	RESULTS	21
7	DISCUSSION	22
8	CONCLUSION AND FUTURE WORKS	24
9	REFERENCES	25

## LIST OF FIGURES

Figure 2.1 DeepSpeech2 Architecture [4].....	6
Figure 2.2 Wav2Vec 2.0 Architecture [5].....	7
Figure 3.1 Overview of low-resource transfer learning .....	9
Figure 3.2 Overview of proposed multi-ASR transfer learning architecture.....	10
Figure 3.3 Architecture of the dual-input CNN model .....	11
Figure 3.4 Overview of combined-input model architecture .....	12
Figure 3.5 DeepSpeech2 Feature Padding .....	13
Figure 3.6 Wav2Vec 2.0 Feature Padding .....	13
Figure 3.7 DeepSpeech2 Feature Transition.....	14
Figure 3.8 Combination of Wav2Vec 2.0 and DeepSpeech2 Features.....	14
Figure 5.1 Process of K-Fold Cross-Validation.....	20

## LIST OF TABLES

Table 4.1 Details of the dataset.....	16
Table 5.1 Details of the experiments.....	18
Table 6.1 Details of the experiment results.....	21
Table 6.2 F1-Score per class .....	21
Table 7.1 Difference in the model transcriptions for a given utterance (DS2 - DeepSpeech2, W2V - Wav2Vec 2.0) .....	22

## LIST OF ABBREVIATIONS AND ACRONYMS

Abbreviation	Description
CNN	Convolution Neural Network
ASR	Automatic Speech Recognition
NLU	Natural Language Understanding
SLU	Spoken Language Understanding
ML	Machine Learning
DS2	DeepSpeech2 Model
W2V	Wav2Vec 2.0 Model
Exp	Experiment
WER	Word Error Rate
LSTM	Long Short-Term Memory
RNN	Recurrent Neural Network
CTC	Connectionist Temporal Classification
SMBO	Sequential Model-based Optimization