

SPRINGERBRIEFS IN SPEECH TECHNOLOGY

STUDIES IN SPEECH SIGNAL PROCESSING, NATURAL
LANGUAGE UNDERSTANDING, AND MACHINE
LEARNING

Manjunath K.E.

Multilingual Phone Recognition in Indian Languages



Springer

SpringerBriefs in Speech Technology

Studies in Speech Signal Processing, Natural Language Understanding, and Machine Learning

Series Editor

Amy Neustein, Fort Lee, NJ, USA

SpringerBriefs present concise summaries of cutting-edge research and practical applications across a wide spectrum of fields. Featuring compact volumes of 50 to 125 pages, the series covers a range of content from professional to academic. Typical topics might include:

- A timely report of state-of-the-art analytical techniques
- A bridge between new research results, as published in journal articles, and a contextual literature review
- A snapshot of a hot or emerging topic
- An in-depth case study or clinical example
- A presentation of core concepts that students must understand in order to make independent contributions

Briefs are characterized by fast, global electronic dissemination, standard publishing contracts, standardized manuscript preparation and formatting guidelines, and expedited production schedules.

The goal of the **SpringerBriefs in Speech Technology** series is to serve as an important reference guide for speech developers, system designers, speech engineers and other professionals in academia, government and the private sector. To accomplish this task, the series will showcase the latest findings in speech technology, ranging from a comparative analysis of contemporary methods of speech parameterization to recent advances in commercial deployment of spoken dialog systems.

More information about this series at <http://www.springer.com/series/10043>

Manjunath K. E.

Multilingual Phone Recognition in Indian Languages

 Springer

Manjunath K. E.
U R Rao Satellite Centre
Indian Space Research Organisation
Old Airport Road, Bengaluru
Karnataka, India

ISSN 2191-737X ISSN 2191-7388 (electronic)
SpringerBriefs in Speech Technology
ISBN 978-3-030-80740-5 ISBN 978-3-030-80741-2 (eBook)
<https://doi.org/10.1007/978-3-030-80741-2>

© The Author(s), under exclusive license to Springer Nature Switzerland AG 2022

This work is subject to copyright. All rights are solely and exclusively licensed by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors, and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

Preface

India is a land of many languages, among them 122 languages are spoken by at least 10,000 people each, with 22 of them constitutionally recognised. Several of the Indian languages do not have sufficient labelled data to develop a separate phone recogniser for themselves. This necessitates an investigation into alternative ways of performing phone recognition, such as multilingual phone recognition. A Multilingual Phone Recognition System (Multi-PRS) is a language-independent, universal Phone Recognition System (PRS) that can recognise the phonetic units present in a speech utterance independent of the language of the speech utterance.

In this book, various aspects of multilingual phone recognition such as development, analysis, performance improvement, and applications of Multi-PRSs are studied for *six* Indian languages – Kannada (KN), Telugu (TE), Bengali (BN), Odia (OD), Urdu (UR), and Assamese (AS). Among the *six* Indian languages considered, Chaps. 3, 4, and 5 use only *four* languages (KN, TE, BN, OD), while the Chap. 6 uses all the *six* languages. The International Phonetic Alphabets (IPA) based transcription is used for deriving a *common multilingual phone-set* by grouping the acoustically similar phonetic units from multiple languages. Both *Gaussian Mixture Model (GMM)-Hidden Markov Models (HMM)* and *Deep Neural Network (DNN)-HMMs* are explored for training the Multi-PRSs using Mel-frequency Cepstral Coefficients (MFCCs) as features under *context independent* and *context dependent* settings. The behaviour of Multi-PRSs across *two* language families namely – Dravidian and Indo-Aryan – is studied and analysed by developing separate Multi-PRSs for Dravidian and Indo-Aryan language families. The performance of Multi-PRSs is analysed and compared with that of the Monolingual Phone Recognition Systems (Mono-PRS).

Articulatory Features (AFs) are explored to improve the performance of Multi-PRSs. The AFs for five AF groups – place, manner, roundness, frontness, and height – are predicted from the MFCCs using DNNs. The oracle AFs, which are derived from the ground truth IPA transcriptions, are used to set the best performance realizable by the predicted AFs. The performance of predicted and oracle AFs are compared. In addition to the AFs, the phone posteriors are explored to further boost the performance of Multi-PRS. Multitask Learning (MTL) is explored to improve

the prediction accuracy of AFs and thereby reducing the Phone Error Rate (PER) of Multi-PRS. Fusion of AFs is done using two approaches: (i) lattice rescoring approach and (ii) AFs as tandem features. It is found that the oracle AFs by feature fusion with MFCCs offer a remarkably low target PER of 10.4%, which is 24.7% absolute reduction compared to baseline Multi-PRS with MFCCs alone. The fusion of phone posteriors and the AFs derived from the MTL yields the best performance. The best performing system using predicted AFs has shown reduction of 3.2% in absolute PER (i.e. 9.1% reduction in relative PER) compared to baseline Multi-PRS.

Applications of multilingual phone recognition in code-switched and non-code-switched scenarios are discussed. Two different approaches for multilingual phone recognition using code-switched and non-code-switched test sets are compared and evaluated. First approach is a front-end Language Identification (LID) system followed by monolingual phone recognisers (LID-Mono) trained individually on each of the languages present in multilingual dataset, while the second approach uses a *common multilingual phone-set* without requiring a front-end LID based switching. Bilingual code-switching experiments are conducted using the code-switched test sets of Kannada and Urdu languages. The state-of-the-art i-vectors are used to perform LID in first approach. It is found that the performance of *common multilingual phone-set* based approach is superior compared to more conventional LID-Mono approach in both non-code-switched and code-switched scenarios. The performance of LID-Mono approach heavily depends on the accuracy of the LID system, and the LID errors cannot be recovered. However, the *common multilingual phone-set* based approach by virtue of not having to do a front-end LID switching and designed based on the common multilingual phone-set derived from several languages is not constrained by the accuracy of the LID system, and hence performs effectively on non-code-switched and code-switched speech, offering low PERs than the LID-Mono system.

This book is mainly intended for researchers working in the area of multilingual speech recognition. This book will be useful for the young researchers who want to pursue research in speech processing with an emphasis on multilingual speech recognition. Hence, this may be recommended as a text or reference book for the postgraduate-level advanced speech processing course. The book has been organised as follows:

Chapter 1 introduces basic concepts of multilingual speech recognition. The multilingual AFs and code-switched speech recognition are briefly introduced. Chapter 2 describes the prior work on multilingual speech recognition systems with primary focus on multilingual AFs and code-switched speech recognition. Chapter 3 describes the development, evaluation, and analysis of Multi-PRS for four Indian languages. Chapter 4 discusses the proposed approaches to derive the multilingual AFs from spectral features. Chapter 5 discusses the use of predicted multilingual AFs to improve the performance of Multi-PRS. Chapter 6 describes the applications of multilingual phone recognition in code-switched and non-code-switched scenarios. Two approaches for multilingual phone recognition are compared using code-switched and non-code-switched test sets. Chapter 7 provides

a brief summary and conclusion of the book with a glimpse towards the scope for possible future work.

I am grateful to my PhD supervisors Prof. V. Ramasubramanian and Prof. Dinesh Babu Jayagopi at the International Institute of Information Technology, Bangalore (IIITB), for their constant support, guidance, and encouragement to carry out this work. This book is based on my doctoral thesis work. I am also grateful to my MS supervisor Prof. K. Sreenivasa Rao at IIT Kharagpur for providing the speech corpora to carry out this work. I thank all the professors and research scholars of IIITB who have helped me to carry out this work. Special thanks to ISRO management and to my colleagues at URSC, ISRO for their cooperation and encouragement during the course of editing and publishing this book. Last but not the least, I am grateful to my parents, my in-laws, my wife, and my daughter for their constant support and encouragement. Finally, I thank all my friends and well-wishers.

Bengaluru, India

Dr. Manjunath K. E.

Contents

1	Introduction	1
1.1	Multilingual Phone Recognition	1
1.2	Articulatory Features for Multilingual Phone Recognition	2
1.3	Approaches for Multilingual Phone Recognition	4
1.4	Code-switched Phone Recognition using Multilingual Phone Recognition Systems	4
1.5	Objective and Scope of the Work	5
1.6	Proposed Organization of the Book	7
	References	8
2	Literature Review	13
2.1	Introduction	13
2.2	Prior Work on Multilingual Speech Recognition	13
2.3	Prior Work on Multilingual Speech Recognition using Articulatory Features	18
2.4	Prior Work on Code-Switched Speech Recognition using Multilingual Speech Recognition Systems	21
2.5	Summary	23
	References	23
3	Development and Analysis of Multilingual Phone Recognition System	27
3.1	Introduction	27
3.2	Experimental Setup	27
3.2.1	Multilingual Speech Corpora	28
3.2.2	Extraction of Mel-frequency Cepstral Coefficients	29
3.2.3	Training HMMs and DNNs	29
3.3	Development of Phone Recognition Systems	30
3.3.1	Development of Monolingual Phone Recognition Systems	30
3.3.2	Development of Multilingual Phone Recognition Systems	32
3.3.3	Development of Tandem Multilingual Phone Recognition Systems	33

3.4	Performance Evaluation of Phone Recognition Systems	34
3.4.1	Performance Evaluation of Monolingual Phone Recognition Systems	35
3.4.2	Performance Evaluation of Multilingual Phone Recognition Systems	35
3.4.3	Performance Evaluation of Tandem Multilingual Phone Recognition Systems	36
3.5	Discussion of Results	37
3.5.1	Analysis and Comparison of the Results	37
3.5.2	Cross-Lingual Analysis	39
3.6	Summary	44
	References	44
4	Prediction of Multilingual Articulatory Features	47
4.1	Introduction	47
4.2	Articulatory Features Specification	47
4.3	Articulatory Feature Predictors (AF-Predictors)	48
4.3.1	Development of Articulatory Feature Predictors	49
4.3.2	Oracle Articulatory Features	51
4.3.3	Performance Evaluation of AF-Predictors	52
4.4	Performance Improvement of AF-Predictors using Multitask Learning (MTL)	54
4.5	Summary	55
	References	55
5	Articulatory Features for Multilingual Phone Recognition	57
5.1	Introduction	57
5.2	Proposed Approaches for Multilingual Phone Recognition using Articulatory Features	57
5.2.1	Development of AF-Based Tandem Multilingual Phone Recognition Systems	59
5.2.2	Fusion of AFs from Multiple AF Groups	61
5.3	Multitask Learning Based AFs for Multilingual Phone Recognition	64
5.4	Summary	65
	References	65
6	Applications of Multilingual Phone Recognition in Code-Switched and Non-code-Switched Scenarios	67
6.1	Introduction	67
6.2	Experimental Setup	67
6.2.1	Multilingual Speech Corpora	68
6.2.2	Code-Switched Test Set	68
6.2.3	Training Support Vector Machines (SVMs)	69
6.2.4	Extraction of i-vectors	70

- 6.3 Approaches for Multilingual Phone Recognition 71
 - 6.3.1 LID-switched Monolingual Phone Recognition (LID-Mono) Approach 72
 - 6.3.2 Multilingual Phone Recognition using Common Multilingual Phone-set (Multi-PRS) Approach 75
- 6.4 Evaluation and Comparison of LID-Mono and Multi-PRS Approaches 76
 - 6.4.1 Non-Code-Switched Scenario 76
 - 6.4.2 Code-Switched Scenario 78
- 6.5 Summary 81
- References 81
- 7 Summary and Conclusion** 85
 - 7.1 Summary of the Book 85
 - 7.2 Contributions of the Book 86
 - 7.3 Future Scope of Work 87
 - Reference 88
- A Support Vector Machines** 89
- B Hidden Markov Models for Speech Recognition** 91
 - Reference 92
- C Deep Neural Networks for Speech Recognition** 93
 - C.1 FeedForward Neural Networks 93
 - C.2 Training Deep Neural Networks 96
 - C.3 Interfacing DNN with HMM (DNN-HMMs) 96
 - References 97
- Index** 99

Acronyms

AF	Articulatory Feature
AF-PER	AF-Prediction Error Rate
AF-Predictor	Articulatory Feature Predictor
AF-Tandem	Combination of AFs as Tandem features
AS	Assamese
ASR	Automatic Speech Recognition
BN	Bengali
CD	Context-Dependent
CI	Context-Independent
DNN	Deep Neural Network
DP	Dynamic Programming
FFNN	Feed-Forward Neural Network
GMM	Gaussian Mixture Model
GMM-UBM	Gaussian Mixture Universal Background Model
HL	Hidden Layer
HMM	Hidden Markov Model
Hz	Hertz
IPA	International Phonetic Alphabet
ISA	Intrinsic Spectral Analysis
KN	Kannada
LFV	Language Feature Vector
LID	Language Identification
LID-Mono	LID-switched Monolingual Approach
LRA	Lattice Rescoring Approach
LVCSR	Large Vocabulary Continuous Speech Recognition
MFCC	Mel-frequency Cepstral Coefficient
MLP	Multi-layer Perceptron
Mono-PRS	Monolingual Phone Recognition System
ms	Millisecond
MSE	Mean Squared Error
MTL	Multitask Learning