

**A PROBABILISTIC APPROACH TO SELECT
UNITS BASED ON ACOUSTIC SIMILARITY
FOR SPEECH SYNTHESIS**

A THESIS

submitted by

ANJANA BABU

for the award of the degree

of

MASTER OF SCIENCE

(by Research)



SCHOOL OF COMPUTING AND ELECTRICAL ENGINEERING

INDIAN INSTITUTE OF TECHNOLOGY MANDI

MARCH 2015

*To my parents,
sister and brother*

DECLARATION

I hereby declare that the entire work embodied in this thesis is the result of investigations carried out by me in the **School of Computing and Electrical Engineering**, Indian Institute of Technology Mandi, under the supervision of **Dr. Anil Kumar Sao**, and that it has not been submitted elsewhere for any degree or diploma. In keeping with the general practice, due acknowledgements have been made wherever the work described is based on finding of other investigators.

Mandi, 175 001

Anjana Babu

THESIS CERTIFICATE

This is to certify that the thesis titled **A Probabilistic Approach to Select Units based on Acoustic Similarity for Speech Synthesis**, submitted by **Anjana Babu**, to the Indian Institute of Technology, Mandi, for the award of the degree of **Master of Science** (by Research), is a bonafide record of the research work done by her under my supervision. The contents of this thesis, in full or in parts, have not been submitted to any other institute or university for the award of any degree or diploma.

Mandi, 175 001

Dr. Anil Kumar Sao
(Guide)

Acknowledgements

I would like to express my sincere thanks to my guide Dr. Anil K Sao for the invaluable guidance and help provided. Also, I would like to extend my gratitude to Prof. Hema A. Murthy for guiding me in my work extensively. Most of the research mentioned in this thesis was done at IIT Madras in collaboration with Raghava Krishnan K in Don Lab. So I would like express my gratitude to Raghava Krishnan K and other members of the lab who constantly aided in the work as well as encouraged wholeheartedly.

Most of the work is carried out as part of the consortium project *Development of Text to Speech Systems for Indian Languages* (11(7)2011-HCC(TDIL)). I would like to thank all the consortium members for sharing the data for research and giving valuable suggestions for research. Also I would like to thank all those who are/were part of the TTS project team in IIT Mandi, for the constant help in my research work. Special thanks to Priyanka Kaushal who helped me with some of the experiments and working hand in hand for many experiments on TTS. Also, I would like to thank research scholars Shejin T., Srimanta Mandal, Pulkit Sharma, Vinayak Abrol, Sarvesh Jayaraman and others for their help in research. I would like to thank research scholars and faculty members in the Multimedia Analytics and Systems group for the weekly research meetings, which have been a great platform to discuss research ideas, expand knowledge. Also, I am happy for the enthusiasm shown by the research scholars and students at IIT Mandi for participating in the evaluation of TTS systems. I am grateful to my Academic Progress Committee members for assessing my progress and keeping me on track. Last but not the least, I would like to thank my dear friends for being with me in times of joy and sorrow.

Anjana Babu

ABSTRACT

Keywords: *TTS, Unit selection speech synthesis, Prosody.*

One of the major challenges in Text-to-Speech Synthesis Systems (TTS) is the incorporation of prosody in the synthesised speech. Various techniques based on linguistic characteristics have been proposed in the literature for improving prosody. However, prosody in TTS is still an open problem and this is addressed in this thesis.

In this work, approaches are proposed for improving the naturalness, intelligibility and prosody. It is performed by selecting the sequence of sound units from a large corpus in such a way that acoustic features are made consistent at (a) segmental level, and (b) supra-segmental level. At segmental level, the differences in acoustic features of sound units are reduced over the entire utterance to improve naturalness and intelligibility. At supra-segmental level, units are selected by ensuring the consistency in the differences in acoustic features of adjacent syllables at phrase level. In this method, consistency of acoustic features is also maintained at utterance level. Unlike the existing USS based TTS, which rely mostly on linguistic information for improving prosody, the proposed approach makes use of acoustic information. Probabilistic approaches are proposed for selecting units based on an acoustic framework. Since the context is only specified by acoustic features, the proposed approaches can be applied to any language and perhaps even for multilingual synthesis. The experimental results of the proposed approaches are demonstrated using five Indian languages. It was observed from the subjective evaluation tests that f_0 contributed to the naturalness of the systems whereas duration and energy helped in improving the intelligibility of the systems. Also, ensuring consistency of energy and f_0 across syllables in phrase and duration across syllables in utterance further improved the prosody.

Contents

Abstract	ii
List of Tables	vi
List of Figures	ix
1 Introduction	1
1.1 Introduction to Text to Speech Synthesis Systems	1
1.2 Prosody in Speech	2
1.3 Objective and Scope of the Work	3
1.4 Contribution of the Thesis	3
1.5 Organisation of the Thesis	5
Abbreviations	1
2 Review of TTS Systems	6
2.1 Introduction	6
2.2 History of speech synthesis	6
2.2.1 Unit Selection Speech Synthesis	7
2.2.2 Statistical Parametric Speech Synthesis	9
2.2.3 Articulatory Synthesis	10
2.3 Prosody Modelling	11
2.3.1 Language based Models	11
2.3.2 Prosodic Event based Models	14
2.4 TTS for Indian Languages	19
2.5 Prosody for TTS in Indian Languages	20

2.5.1	Prosodic Phrases and Pause Duration	21
2.5.2	Duration	23
2.5.3	f_0 Pattern	24
2.5.4	Energy Modification	24
2.6	Speech Database	24
2.7	Summary	25
3	Probabilistic Approach to Select Units Based on Acoustic Similarity	27
3.1	Introduction	27
3.2	Potential Causes for Lack of Naturalness and Intelligibility	28
3.2.1	Duration of Syllables in Utterance	29
3.2.2	f_0 of Syllables in Utterance	31
3.2.3	Energy of Syllables in Utterance	38
3.3	Proposed Approach to Select Units	43
3.4	Important Aspects of Implementation	46
3.4.1	Distribution of Acoustic Features	46
3.4.2	Context information of Syllables	51
3.4.3	Selecting Units for Synthesis	51
3.4.4	Experimental Results	53
3.5	Spectral Continuity	57
3.5.1	Results	58
3.6	Summary	60
4	Probabilistic Approach to Select Units using Phrase Information	62
4.1	Introduction	62
4.2	Variation of Acoustic Features at Supra- segmental Level	63
4.2.1	f_0 of Phrases	64
4.2.2	Energy of Phrases	64
4.2.3	Duration of Phrases	65
4.3	Proposed Approach	66
4.4	Evaluation	71
4.5	Summary	73

5 Summary and Conclusion	74
5.1 Summary	74
5.2 Future Works	75
References	88
List of Publications	89