*Department of Electrical & Computer Engineering*
*The University of Auckland*
*New Zealand*

# Embedded Speech Recognition Systems

*Octavian Cheng*

*September 2008*

*Supervisors:   Dr Waleed Abdulla*

*Prof Zoran Salcic*

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS OF DOCTOR OF PHILOSOPHY IN ENGINEERING

# Abstract

Despite many research efforts, automatic speech recognition (ASR) is still not widely used in embedded systems. One of the reasons for this is that embedded platforms are limited in terms of processing power and computing resources. These limitations contribute to an increase in decoding time, which can lead to poor user experiences. As a result, a compromise is often made between decoding speed and other performance criteria, such as recognition accuracy and vocabulary size.

There are two main objectives in this thesis. The first objective is to develop an embedded ASR system which is suitable for real-time applications. We focus on various kinds of approaches that can reduce the decoding time without severe performance degradation in recognition accuracy and vocabulary size. The task is a 993-word Resource Management (RM1) task, which serves as the benchmark for command-and-control type of applications. The target platform is an Altera Nios II softcore processor system running at 120MHz. The system is synthesized on a Stratix II FPGA device.

Three approaches have been successfully adopted on the target platform. First, due to lack of hardware support for floating-point arithmetic, we propose a framework for converting various data types to fixed-point formats. Experimental results on the RM1 task show that the fixed-point system is 9 times faster than the floating-point system without any degradation in recognition accuracy. Second, a hardware-software co-processing ASR system is developed. The architecture mainly consists of a Nios II processor and a hardware accelerator. The hardware accelerator is responsible for the calculation of the Gaussian mixture model (GMM) emission probabilities, which is the major computational bottleneck. The co-processing system is tested on the same RM1 task. In comparison with the pure software-based fixed-point system, the average real-time factor improves from 1.87 to 0.62 (about 3 times speed-up). The word accuracy rate is 93.33%, which is the same as that of the pure software-based system. Third, in order to further improve the timing performance, an adaptive beam pruning algorithm is introduced, which applies tighter pruning when there are a large number of active hypotheses. For the same RM1 task, the average real-time factor further reduces to 0.54. The word accuracy rate is 93.16%.

Apart from recognition accuracy, decoding speed and vocabulary size, another point of consideration when developing a practical ASR application is the adaptability of the system. An ASR system is more useful if it can cope with changes that are introduced by users, for example, new words and new grammar rules. In addition, the system can also automatically update the underlying knowledge sources, such as language model probabilities, for better recognition accuracy. Since the knowledge sources need to be adaptable, it is inflexible to statically combine them. It is because on-line modification becomes difficult once all the knowledge sources have been combined into one static search space.

The second objective of the thesis is to develop an algorithm which allows dynamic integration of knowledge sources during decoding. In this approach, each knowledge source is represented by a weighted finite state transducer (WFST). The knowledge source that is subject to adaptation is factorized from the entire search space. The adapted knowledge source is then combined with the others during decoding. In this thesis, we propose a generalized dynamic WFST composition algorithm, which avoids the creation of non-coaccessible paths, performs weight look-ahead and does not impose any constraints to the topology of the WFSTs. Experimental results on Wall Street Journal (WSJ1) 20k-word trigram task show that our proposed approach has a better word accuracy versus real-time factor characteristics than other dynamic composition approaches.

# Acknowledgements

I would like to thank my thesis supervisors, Dr Waleed Abdulla and Prof Zoran Salcic, for their guidance and support. Throughout my PhD years, I learnt a lot from them. I always received many useful suggestions from them in our meetings, especially when problems occurred and research results were not as promising as expected.

I also would like to thank Prof Hervé Bourlard for offering me a research internship at IDIAP, Switzerland. The stay at IDIAP broadened my view on speech recognition research. I would like to acknowledge the members of IDIAP. In particular, I would like to thank John Dines and Mathew Magimai Doss for the collaboration.

I would like to thank New Zealand Tertiary Education Commission for the scholarship offer. With their financial support, I can focus solely on my thesis work.

Finally, I would like to thank my parents for their support and encouragement.

# Contents

# List of Figures

# List of Tables

# Glossary

| | |
|---|---|
| ALM | Adaptive logic module |
| ANN | Artificial neural network |
| API | Application programming interface |
| ASR | Automatic speech recognition |
| CD | context-dependent |
| CI | context-independent |
| CMN | Cepstral mean normalization |
| CMVN | Cepstral mean variance normalization |
| DTW | Dynamic time warping |
| EM | Expectation-maximization |
| FPGA | Field programmable gate arrays |
| GMM | Gaussian mixture model |
| HDL | Hardware description language |
| HMM | Hidden Markov model |
| IP | Intellectual property |
| LM | Language modelling |
| LPC | Linear predictive coding |
| MFCC | Mel frequency cepstral coefficient |
| MLP | Multi-layer perceptrons |
| PLP | Perceptual linear prediction |
| RTF | Real-time factor |
| SOPC | System-on-a-programmable-chip |
| VHDL | VHSIC hardware description language |
| VHSIC | Very-high-speed integrated circuits |
| WER | Word error rate |
| WFST | Weighted finite state transducer |