# Cairo University
## Institute of Statistical Studies and Research
## Department of Computer and Information Sciences

# BUILDING A SPEECH RECOGNITION SYSTEM FOR SPOKEN ARABIC

**by**

**Tarek Ahmed Fouad Ibrahim Sheisha**

Under the Supervision of

**Prof. Atef  M. A-Moneim**      **Prof. Khaled Fouad Shaalan**

Institute of Statistical Studies and Research    Faculty of Computers and Information
Cairo University            Cairo University

A Thesis Submitted to the
Department of Computer and Information Sciences
In Partial Fulfillment of the
Requirements for the Degree of
MASTER OF SCIENCE
In
Computer Science

Cairo
**August 2009**

# Approval Sheet

## BUILDING A SPEECH RECOGNITION SYSTEM FOR SPOKEN ARABIC

By
Tarek Ahmed Fouad Ibrahim Sheisha

A Thesis Submitted to the Department of Computer and Information Sciences, Institute of Statistical Studies and Research, Cairo University in Partial fulfillment of the Requirements for the Degree of Master of Science in Computer Science.

**Approved by The Examining Committee:**

| Name | Signature |
|------|-----------|
| Prof.  Atef M. A-Moneim | |
| Dr.  Hesham N. Elmahdy | |
| Dr.  Hesham A. Hefny | |

Cairo
August, 2009

# **<u>Statement</u>**

I certify that this work has not been accepted in substance for any academic degree and is not being concurrently submitted in candidature for any other degree.

Student Name:     Tarek Ahmed Fouad Ibrahim Sheisha

Signature:     *Tarek A. Fouad*

# **<u>Acknowledgements</u>**

# ABSTARCT

The automatic recognition of spontaneously spoken speech is a very challenging task because it contains much more disfluencies than the read speech. The difficulty is even increased in case of recognizing colloquial speech. The colloquial speech contains many changes from the standard language such as changes in the syntax, the structure, the letters, and the diacritics. Moreover, the colloquial speech may contain different expressions and slang. Building of Arabic speech system is useful for adding user-friendly interfaces for Arabic software. Also, the embedding of the Arabic speech recognition system in the machine translation systems can allow human-to-human communication across languages boundaries.

In this thesis, a system was built for recognizing spontaneous Egyptian colloquial Arabic. A speech database was collected by recording 20 hours of spontaneously spoken Arabic dialogues. All spoken dialogues were segmented into separate utterances, the corresponding transcriptions were written. The diacritics (FatHa, Dammah, and Kasrah) are included in the transcription of each utterance. The speech database was used to train the acoustic models and the language model, which were used as knowledge sources for building the system. A baseline dictionary was built. A set of phones was chosen, which is suitable for the phonetics of Egyptian colloquial Arabic. A statistical language model, a word bigram Language model, was built using the canonical-form transcriptions of utterances of the speech database. A statistical acoustic model was built for each phone. The acoustic models were based on the Hidden Markov Models (HMMs). There are separate HMMs for the diacritics. All the acoustic models were trained using the speech database. Another system was built through training the models with non-diacritized transcriptions (i.e. the short vowels were removed).

The recognition experiments revealed that the system which was built with the fully diacritized transcriptions had a higher percentage of correctly recognized words. The system was modified by converting the monophone units into triphones units. Another modification was to increase the number of Gaussian

mixtures associated with each HMM's state. The optimum values of the word insertion penalty P, and the grammar scale factor S were experimentally found. The percentage of correctly recognized words for the baseline system is 82.87 %.

A proposed data-driven method of modeling pronunciation variation is utilized to improve the baseline system. This method modifies the three levels of the speech recognition system at which modeling can take place; i.e. the dictionary, the acoustic model, and the language model.

 A dynamic programming-based algorithm was developed to align the canonical-form transcription of each utterance with the string of phones constituting that utterance. For every utterance, a string of the constituting phones was obtained by a developed phone recognizer. As a result of the alignment process, a set of pronunciation rules was generated. The set of rules was applied on the whole set of canonical-form transcriptions to generate the surface-form transcriptions. The acoustic and language models were re-trained with the surface-form transcriptions. Then, by applying suitable rules on each entry in the dictionary, some pronunciation variants are generated. The dictionary was also modified by adding the pronunciation variants to the baseline dictionary.

The recognition experiments showed that the modeling of the pronunciation variations improved the percentage of correctly recognized words to 91.9 %.

A user interface was developed, through which the recognized sentences are displayed in diacritized Arabic scripts.

# Contents

**Chapter 3 Search Strategies in Speech Recognition Systems**

**Chapter 4 Arabic Speech**

**Chapter 7 Conclusion and Future Work**

# List of Figures

# List of Tables

# List of Algorithms