



**Information Systems Department  
Faculty of Computer and Information Sciences  
Ain Shams University**

# **Content Based Search Engine for Arabic Audio Files**

Thesis submitted as a partial fulfillment of the requirements for the  
degree of Master of Science in Computer and Information Sciences

By

**Mona A.Azim A.Gawad**

Teaching Assistant at Information Systems Department,  
Faculty of Computer and Information Sciences,  
Ain Shams University

Under Supervision of

**Prof. Mohamed Fahmy Tolba**

Scientific Computing Department,  
Faculty of Computer and Information Sciences,  
Ain Shams University

**Prof. Tarek Fouad Gharib**

Information Systems Department  
and Head of Information Systems Department,  
Faculty of Computer and Information Sciences,  
Ain Shams University

**Prof. Nagwa Lotfy Badr**

Information Systems Department  
and Vice Dean of education and student affairs,  
Faculty of Computer and Information Sciences,  
Ain Shams University

April – 2017  
Cairo

# Acknowledgment

First and foremost, my deep gratitude to Allah, the most beneficial, the most graceful, for helping me completing this work to the fullest.

Second, I would like to express my veracious and deep respect and gratitude to all my supervisors; *Prof. Mohamed Fahmy Tolba*, *Prof. Tarek Fouad Gharib* and *Prof. Nagwa Lotfy Badr* at the Faculty of Computer and Information Sciences, Ain Shams University, who dedicated much of their precious time to assist me in elaborating this thesis, guide me in understanding and analyzing problems, and continuously provide support valuable comments that influenced my career and will be my success recipe throughout my entire life.

Finally, I would like to acknowledge Dr. Abd El Aziz Abd El Hamid at the Faculty of Computer and Information Sciences, Ain Shams University for his continuous support and efforts throughout this work until it reached the desired form. In addition, I would like to acknowledge my dearly loved parents who were burdened with even more roles and responsibilities in order to give me the chance to accomplish this work and now I hope that they are proud to see this work has been finished successfully.

# **Abstract**

Information systems were subject to considerable changes. These changes were arisen due to the fact that information representation underwent massive evolution. This evolution from heavyweight information representation (i.e. papers, books, music sheets) towards lightweight information (i.e. multimedia and digital forms) is driven by spectacular changes to the constraints of storing, managing and retrieving such types of information.

Multimedia Information Retrieval (MIR) is about obtaining semantic information from different multimedia data sources (i.e. audio, video or images). The research in this field faces many challenges due to the variety of data types, methodologies and research problems. One of the challenging sub-field of MIR is audio retrieval.

Audio retrieval aims to extract certain piece of information from audio signals contents based on a submitted query. Audio data types may include pure music, environmental sounds, songs and speech utterances.

In this work, we have implemented a content based retrieval system for Arabic speech signals to retrieve the speech audios relevant to the user text query. A Large Vocabulary Arabic Continuous Speech Recognition (LVACSR) System was built as a preceding step towards defining speech transcriptions used in the retrieval system. Hidden Markov Models (HMMs) based LVACSR system has adopted a novel approach in state tying that uses a new developed phonetic decision tree for the Arabic language for building acoustic models of the LVACSR system. Building such models was accomplished by firstly building phoneme-based HMMs then building triphones based HMMs. The development of the triphones models process includes using both Data driven based and phonetic tree based state tying techniques. After testing the built models, the tree-based

tied states models have shown maximum word correctness and accuracy. The text transcriptions generated from the recognition system that used the tree-based tied states models as acoustic models were used as an input to the text-based retrieval system. In the text-based retrieval system, the text transcriptions were stemmed then indexed in the indexer database then the most relevant audios to the user query were retrieved. The text-based retrieval system used three different character based similarity distance to obtain the desired audios list. Levenshtein distance, Jaro-Wrinkel distance and N-gram distance were calculated to measure the similarity between the user query and the indexed text transcriptions of the audios. After obtaining the most relevant audios list, the precision and recall were calculated to evaluate the text-based retrieval system. The results showed that the Levenshtein distance and Jaro-Wrinkel distance were more accurate than the N-gram distance when used as a similarity measure.

## List of Publications

Azim, M. A., Hamid, A. A. A., Badr, N. L., & Tolba, M. F. (2016, December). Large Vocabulary Arabic Continuous Speech Recognition Using Tied States Acoustic Models. *Asian Journal of Information Technology*.(Accepted on December 10,2016)

Azim, M. A., Hamid, A. A. A., Badr, N. L., & Tolba, M. F. (2016, October). Tree-Based HMM State Tying for Arabic Continuous Speech Recognition. In *International Conference on Advanced Intelligent Systems and Informatics* (pp. 96-103). Springer International Publishing.

Azim, M. A., Badr, N. L., & Tolba, M. F. (2016, May). An Enhanced Arabic Phonemes Classification Approach. In *Proceedings of the 10th International Conference on Informatics and Systems* (pp. 210-214). ACM.

# Table of Contents

Acknowledgment .....	II
Abstract .....	III
List of Publications .....	V
Table of Contents .....	VI
List of Figures .....	VIII
List of Tables .....	IX
List of Abbreviations .....	X
Chapter 1. Introduction .....	2
1.1 Overview .....	2
1.2 Motivation .....	2
1.3 Objectives.....	3
1.4 Thesis Organization .....	3
Chapter 2. Literature Review.....	5
2.1 Content based audio retrieval systems .....	5
2.2 Arabic speech recognition systems .....	6
2.3 State tying techniques .....	8
Chapter 3. The Proposed Content Based Arabic Speech Retrieval System.....	11
3.1 System Architecture .....	11
3.2 LVACSR System .....	13
3.2.1 Feature Extraction:.....	14
3.2.2 Training Acoustic Models: .....	16
3.2.2.1Acoustic Models.....	16
3.2.2.2HMMs definition .....	16
3.2.2.3 HMMs parameters estimation.....	18
3.2.2.4HMM State tying Techniques.....	20
3.2.3 Building Phonetic Dictionary.....	23

3.2.4	Building Language Model .....	25
3.2.5	Recognizer: .....	26
3.2.5.1	HMMViterbi Decoding:.....	26
3.3	Text based Retrieval System.....	27
3.3.1	Stemmer and Indexer .....	27
3.3.2	Semantic Similarity Retrieval .....	28
3.3.2.1	Levenshtein Distance Similarity: .....	28
3.3.2.2	Jaro Distance: .....	29
3.3.2.3	Jaro–Winkler Distance: .....	29
3.3.2.4	N-gram Similarity .....	29
Chapter 4.	Experimental Results .....	31
4.1	Data set.....	31
4.2	HTK Toolkit.....	33
4.3	LVACSR system Evaluation .....	34
4.4	Text based Retrieval System Evaluation .....	42
4.5	Results and discussion.....	42
Chapter 5.	Conclusion and Future Work .....	48
Chapter 6.	References .....	53
APPENDIX A	Arabic Phonetic Decision Tree based Tying Commands	59

## List of Figures

Figure 3.1 The proposed system architecture.....	12
Figure 3.2 A sample run of the text based retrieval system for the query "خدمة".....	13
Figure 3.3 A typical 3-state right to left HMM.....	17
Figure 3.4 A sample of the phonetic decision tree.....	23
Figure 3.5 Tree-Based state tying in HMMs.....	25
Figure 3.6 A sample of the phonetic dictionary items.....	25
Figure 4.1 A sample from the text corpus.....	33
Figure 4.2 Phoneme based transcription for the word "الْمُهَيَّئَةُ".....	33
Figure 4.3 HMM configuration file.....	35
Figure 4.4 HMM prototype.....	36
Figure 4.5 Tying Transition matrices.....	38
Figure 4.6 Data driven state tying.....	39
Figure 4.7 Word Correctness for the built models.....	44
Figure 4.8 Monophonic based models vs Triphonic based models word Correctness(%).....	45
Figure 4.9 various phoneme sets that corresponds to the word "الشركة".....	46



## **List of Tables**

Table 4.1 General statistics of both text and audio corpora in each category .....	32
Table 4.2 Phonemes and their corresponding labels in the transcription files .....	32
Table 4.3 Arabic Vowels Phonemes Inventory .....	41
Table 4.4 Arabic Consonants Phonemes Inventory .....	41
Table 4.5 Monophonic, data driven based and tree based tied triphones Experiments Results. ....	43
Table 4.6 Sample of the recognizer results using tree based tied states tri-phones .....	44
Table 4.7 Precision and Recall of Text based retrieval Experiments .....	46

## List of Abbreviations

CBASR **C**ontent **b**ased **A**rabic **S**peech **R**etrieval

ASR **A**utomatic **S**peech **R**ecognition.

CL **C**lone **C**ommand

HERest **H**TK **E**mbedded **re-est**imation

HHed **H**MM **H**Tk **E**ditor

HLEd **H**MM **L**abel **E**ditor

HMMs **H**idden **M**arkov **M**odels

HTK **H**MMs **T**ool **K**it

IVR **I**nteractive **V**oice **R**esponse

LVACSR **S**ystem **L**arge **V**ocabulary **A**rabic **C**ontinuous **S**peech  
**R**ecognition **S**ystem

QS **Q**uestion **C**ommand

SASSC **S**tandard **A**rabic **S**ingle **S**peaker **C**orpus

TC **T**riphones **C**onversion

TI **T**ie **C**ommand

WER **W**ord **E**rror **R**ate

LM **L**anguage **M**odel

MFCC **M**el **F**requency **C**epstral **C**oefficients

# **Chapter 1**

---

## **Introduction**

---

# **Chapter 1. Introduction**

## **1.1 Overview**

Nowadays, digital archives are increasing rapidly and the need for effective retrieval systems is increasing correspondingly. In case of audio, traditional audio retrieval systems rely mainly on text-based approaches. These approaches depend on gathering and indexing audio files meta data then applying a text based retrieval technique on the indexed data. However, this approach is simple and straightforward and reuse the already refined text based search techniques, the misplaced audio meta data may result in poor retrieval quality.

Recently, content-based approaches were developed to improve the audio retrieval systems performance. These approaches aim to fulfill the user audio information needs by searching the entire audio library and retrieve the most relevant ones to the user query.

## **1.2 Motivation**

The challenge of audio information retrieval systems is familiar to anyone who has been in a vacation and then after returning home he found his answering machine was full of messages and he has to be patient enough to listen to the entire tape till he found the urgent message from his boss.

Content-based audio retrieval has not been explored in the context of content based multimedia retrieval systems. However, if fully studied it can also be more useful in many applications such as audio data set storing and managing systems and digital libraries manipulation.

Searching and archiving digital libraries is an extensive and exhausting task to be done manually. Thus automating such task is

mandatory for certain beneficial organizations such as Radio stations and recordings incorporations.

### **1.3 Objectives**

This work aims to build a content based Arabic speech retrieval system. The text transcriptions of the Arabic speech audio files were generated using a Large Vocabulary Arabic Continuous Speech Recognition (LVACSR) system. Phonemic based Hidden Markov Models (HMMs) were trained and tested to serve as the LVACSR system acoustic models. In order to improve the system word correctness, triphones based HMMs were built as well. Firstly, the triphones models applied the data driven based state tying. After that, a novel Arabic phonetic decision tree was used to apply the decision tree based state tying on the triphone models. Finally, for each of the trained models the retrieval results were evaluated.

### **1.4 Thesis Organization**

The thesis is organized as follows. Chapter 2 reviews the recent advances in the audio retrieval systems and gives a brief overview of the latest trends of Arabic speech recognition systems. Chapter 3 introduces the proposed content based Arabic audio retrieval system and the typical architecture of the LVACSR with a detailed explanation of each system module and the text based retrieval system. Experimental Results and system evaluation are discussed in chapter 4. Conclusions and future work are presented in chapter 5.

## **Chapter 2**

---

# **Literature Review**

---

## **Chapter 2. Literature Review**

### **2.1 Content based audio retrieval systems**

Most of the information retrieval systems are used to focus on the textual type of information and text based information retrieval techniques have been developed and achieved considerable results.

Due to the increase of multimedia databases, the need of developing efficient retrieval techniques has arisen. Recently, the multimedia information retrieval has gained much interest by various researchers. The multimedia information retrieval is the research field that aims at extracting semantic information from multimedia data sources; such as images, audios and videos. Most of researchers in multimedia information retrieval systems focus on searching the visual media (i.e. image and videos).

However, audio is an important medium and an information carrier from the human auditory view point, so it is needed to search and retrieve the audio collections. In fact, audio can save large amount of rich information in the form of speech, musical and sound effects within a sort of small space. So it should be searched based on its aural content; such as speech uttered within audio, musical melodies and acoustic features.

The development of content based audio retrieval systems depends on the type of the audio either it is a piece of music, an uttered speech or a sound effect or the query type either it is query by example or text query.

The research in this field was initially proposed in [1] by implementing a general audio classification and retrieval system. In that system, audio is represented by perceptual and acoustical features, in which users can search or retrieve sounds by several types of query. In [2] a newly developed