**Ain Shams University**

Faculty of Engineering

Electronics and commuincation Engineering Department

# Synthesis of Arabic Speech Signals

A Thesis submitted in partial fulfillment

of the requirements of the M.Sc. degree in the Electrical Engineering

(Communications)

By

**Eng. Sammar Mohamed Elsayed Soliman**

Eng. General Organization of Educational Buildings

**Supervised by:**
**Prof. Dr. Salwa H. El-Ramly**
Ain Shams University

**Dr. Nemat Sayed Abd El-Kader**
Cairo University

Cairo-1997

## ACKNOWLEDGMENT

The author wishes to acknowledge the sincere effort and continual guidance of her supervisor **Prof. Dr.: Salwa Hussein El-Ramly**. Without her scientific assistance and precious suggestions the thesis would not have that reach

The author also wishes to express her sincere gratitude to **Dr. : Nemat Sayed Abd El-kader** for helping in deciding the subject of the research, defining the problem, useful suggestions, technical support, and continual follow-up during the research

The author wishes to thank the **Department of Electronics and Communications Engineering, Faculty of Engineering, Ain Shams University.**

Furthermore, the author wishes to thank the *General Organization* of **Educational Buildings.Nasr City.**

بسم الله الرحمن الرحيم

## Statement

This dissertation is submitted to Ain Shams University for the M.Sc. degree in Electrical Engineering (Communications Engineering).

The work included in this thesis was carried out by the author at the Electronics and Communications Engineering Department. Ain Shams University.

No part of this thesis has been submitted for a degree or qualification at any other university or institute
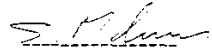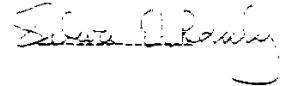
Date:
Signature
Name. Sammar Mohamed Elsayed Soliman

## EXAMINERS COMMITTEE

**Name. Title & Affiliation:**                    Signature

1- Prof. Dr. Safwat Mahrous Mahmoud
   Electronics & communication Eng.Dept.
   Ain Shams University.

2- Prof. Dr. Salwa Hussein El- Ramly
   Electronics & communication Eng.Dept.
   Ain Shams University.

3- Prof. Dr. Mohamed Yones Abd El-Samee Elhamalawy --------
   Computers and system Engineering Dept.
   Elazhar University.

Date :    1997

# ABSTRACT

**Sammar Mohamed Elsayed Soliman. Synthesis of Arabic Speech Signals. Master Degree.** Electronics and Communications Engineering Department. Faculty of Engineering. Ain Shams University. 1997.

Using computers in many fields of our daily life requires direct and easy method for dealing with it. Direct communication with computers by text or speech is still the greatest aim of researchers in this field

During the last twenty years, researchers have developed many techniques to produce computerised speech, requiring relatively small data storage requirements and offering excellent capability for concatenation of phrases

Electronic speech synthesis aims to cover concisely but completely all fundamental subject matter which is useful either for research or for the design of speaking systems

Although many text-to-speech systems are commerically available today for a number of languages. (English, French, German, Japanese, etc.) in Arabic language few research works have been done. Recently, there is a growing interest among computer scientists to develop Arabic text-to speech systems and to synthesize naturally sounding Arabic.

This thesis aims to build a software formant synthesizer that is suitable for the nature of the Arabic language (in spite of its simulation complexity) and using this synthesizer to construct an inventory of the basic Arabic speech units with all phonetic variants of these synthesis units that can be used as a first step of an Arabic text-to-speech system. A comparative study between different basic units is made and then the allophones are chosen as the basic units.

The implemented synthesizer consists of two main steps, extraction of the analysis parameters for each allophone (analysis phase), and the speech reproduction (synthesis phase).
In the analysis phase each allophone is extracted from a set of natural Arabic words containing this allophone in the same place or in various places and then analysed to extract its features (pitch, first four formants, first four bandwidths, voiced/unvoiced classification). A new set of parameters will be stored every 20 ms on the PC hard disk.
Storing these parameters instead of the speech signal itself reduces the required memory storage for the constructed data base.
In the synthesis phase, the files which contain the desired parameters of a new word or utterance, are called and concatenated in a new file that activates the synthesizer to produce the output speech.

Experimental results indicate good output speech quality due to good implementation of the synthesizer and careful choice of the phonetic units. Using the constructed data base with a suitable interpolation method the output speech will be improved.

# Contents

# List of Figures

# List of Tables