

## THESIS EVALUATION FORM

We certify that we have read this thesis and that in our opinion it is fully adequate, in scope and qualify as an undergraduate thesis, based on the result of the oral examination taken place on .../.../2011.

.....

(Advisor)

.....

(Committee Member)

.....

(Committee Member)

Prof. Dr. Gülay TOHUMOĞLU

(Chairman)

## **ACKNOWLEDGEMENTS**

First of all i am glad to my advisor Assist. Prof. Dr. Gulden KÖKTÜRK encouraging me for this project, then i am glad to the committee members Assist. Prof. Dr. Hakkı Tarkan YALAZAN and Haldun SARNEL for spending time and guiding, and then i am glad to my dear friends İrem ÇAMAŞIRCIOĞLU and Kutlu TACER for helping me in the all steps of my project.

**FATİH KARTAL**

## **ABSTRACT**

Turkish songs processing with the speech recognition method is the goal of this project. Speech recognition was a curiosity whether it is possible or not to be realized but nowadays it is used in wide range applications from safety systems to funny games. By speech recognition success, it is desired to construct a “karaoke” system. Project included both hardware and software realization of the system but SONY ® did not reply to my e-mails which are about this project so it was impossible to realize hardware part. This project includes software realization.

Karaoke is a form of interactive entertainment or video game in which amateur singers sing along with recorded music using a microphone and public address system. The music is typically a well-known pop song. Lyrics are displayed on a video screen, along with a moving symbol or changing color to guide the singer. Karaoke is a very popular game between young people. The entertainment industry is very important and the karaoke charms the boss with its portion in the market. In this project, a karaoke system is constructed and as different the singer gets a point at the end of the song. The score is given according to the performance of singer and this performance is evaluated by speech recognition algorithm.

There are several developed techniques in digital communication for Automatic Speech Recognition (ASR). These techniques are developed by using different features of human voice. The speech is a continuous time signal and it is necessary to digitalize it with an appropriate sampling rate in order to be processed by a digital processor. There are lots of factors must be taken into account to develop ASR systems such as vocabulary size, if the system speaker dependent or speaker independent, if the ASR is required to recognize isolated words or continuous speech, and so forth.

In this project, the goal is to realize an ASR system which is speaker independent and is required to recognize isolated words. Word numbers is an important factor and the system performance decreases when the vocabulary size increases. And also the sickness, possible stress of the human, and noise interference affect the system performance.

Lots of factors affect performance of the system so it is hard to construct a perfect speech recognition system. It is important to train all factors and analyze the behavior of system against these factors. The algorithm should give the optimum result under all of these undesired conditions.

## ÖZET

Bu projenin amacı Türkçe şarkıların ses tanıma yöntemi ile işlenmesidir. Ses tanıma geçmişte mümkün olup olmadığı merak edilen bir konuydu ancak günümüzde güvenlik sistemlerinden eğlenceli oyunlara kadar birçok uygulamada kullanılan bir yöntem olmuştur. Ses tanımanın gerçekleştirilmesi yardımıyla bir “karaoke” sistemi gerçekleştirilmek istenmiştir. Projenin amacı bu sistemi hem donanımsal hem de yazılımsal olarak gerçekleştirmektir ancak SONY ®, PS3 oyun konsolunda şarkıların nasıl seçildiği konusunda göndermiş olduğum e-postalara cevap vermediği için bu proje sadece yazılımsal olarak gerçekleştirilmiştir.

“Karaoke” amatör şarkıcıların ekranda gördükleri şarkı sözlerini söyleyerek eğlendikleri bir oyundur. Şarkı genelde yaygın olarak bilinen pop müzik türünden seçilir. Şarkı sözleri ekranda gösterilir ve şarkının zamanlamasına uygun olarak kelimeler değişerek şarkıcıya şarkıyı nasıl söylemesi gerektiği konusunda rehberlik eder. Eğlence sektörünün piyasada büyük bir payı olduğu günümüzde, gençler arasında çok yaygın olan bu oyun patronlarında dikkatini çekmektedir. Bu projede “karaoke” oyununa ek olarak oyunun sonunda şarkıcının performansına göre puan verilmektedir. Bu performansın başarısına projede yapmış olduğum ses tanıma sistemi karar vermektedir.

İnsan sesinin farklı özellikleri kullanılarak otomatik ses tanıma için birçok farklı teknik geliştirilmiştir. Ses sürekli bir sinyaldir ve dijital bir işlemci tarafından işlenebilmesi için uygun örnekleme frekansıyla dijitalleştirilmelidir. Otomatik ses tanıma sistemlerinde kelime sayısı, sistemin konuşmacıya bağlı olup olmadığı, ne tür konuşmaların anlaşılacak istendiği gibi dikkat edilmesi gereken birçok özellik vardır.

Bu projede amaç konuşmacıdan bağımsız ve yalıtılmış kelimeleri tanıyan otomatik ses tanıma sistemini gerçekleştirmektir. Kelime sayısı sistem performansını etkileyen en önemli faktörlerden bir tanesidir ve kelime sayısı ile başarı oranı ters orantılıdır. Ayrıca hastalık, insanın stresli olması, mikrofona olan uzaklık ve ortamdaki gürültüde sistemin performansını olumsuz etkileyen faktörlerdendir.

Sistemin performansını etkileyen birçok faktör olduğu için hatasız bir ses tanıma sistemi gerçekleştirmek oldukça zordur. Bu yüzden sistem bütün bu faktörlerle ayrı ayrı denenmeli ve herbirine verdiği tepki incelenerek bunlar giderilmeye çalışılmalıdır. Önemli olan tüm istenmeyen koşullar altında en iyi performansı verebilecek olan sistemi gerçekleştirmektir.

# TABLE OF CONTENTS

	<b>Page</b>
Thesis Evaluation Form -----	<b>i</b>
Acknowledgements-----	<b>ii</b>
Abstract-----	<b>iii</b>
Özet-----	<b>iv</b>
Table of Contents-----	<b>v</b>
List of Figures-----	<b>viii</b>
List of Tables-----	<b>IX</b>
 <b>1.INTRODUCTION</b>	
1.1 Inspiration of the Project -----	1
1.2 General Overview of the Project -----	1
1.3 Aim of the Project -----	3
 <b>2. METHODOLOGY</b>	
2.1 Production of the Speech -----	4
2.2 Technical Characteristics of the Speech Signal -----	5
2.2.1 Bandwidth-----	6
2.2.2 Fundamental Frequency-----	6
2.2.3 Peaks in the Spectrum-----	6
2.2.4 The Envelope of the Power Spectrum Decreases with Increasing Frequency--	6
2.3 Pre-processing-----	7
2.4 Feature Extraction-----	7

2.4.1 Windowing-----	7
2.4.2 Spectrum Computation (using FFT)-----	13
2.4.3 Power Spectrum Computation-----	14
2.5 Features and Vector Spaces-----	14
2.5.1 Feature Vectors and Vector Spaces-----	14
2.5.2 A Simple Classification Problem-----	15
2.6 Classification of Vectors-----	16
2.6.1 Prototype Vectors-----	16
2.6.2 Nearest Neighbor Classification-----	16
2.7 Distance Measures-----	17
2.7.1 Euclidean Distance-----	17
2.7.2 City Block Distance-----	17
2.7.3 Weighted Euclidean Distance-----	18
2.7.4 Mahalanobis Distance-----	19
2.8 Dynamic Time Warping-----	24
2.8.1 Dynamic Time Warping (DTW)-----	26
2.8.2 DTW Algorithm-----	28
2.8.3 Complexity of DTW-----	32
2.8.4 Speeding up DTW-----	32
2.8.5 Fast DTW Algorithm-----	34
<b>3. RESULTS-----</b>	<b>37</b>
<b>4. COST ANALYSIS-----</b>	<b>38</b>

<b>5. CONCLUSIONS</b>	<b>39</b>
<b>6. REFERENCES</b>	<b>39</b>

## List of Figures

- Figure1:** A simple block diagram of the speech recognition system (page 4)
- Figure2.1:** A schematic diagram of the human speech production apparatus (page 5)
- Figure2.3:** The hamming window (page 9)
- Figure2.4:** Rectangular windowed sine wave (page 10)
- Figure2.5:** Hamming windowed sine wave (page 10)
- Figure2.6:** FFT of rectangular windowed sine wave (page 11)
- Figure2.7:** FFT of hamming windowed sine wave (page 11)
- Figure2.8:** FFT of rectangular windowed sine wave in dB (page 12)
- Figure2.9:** FFT of hamming windowed sine wave in dB (page 12)
- Figure2.10:** A map of feature vectors (page 15)
- Figure2.11:** Two dimensions with different scales (page 18)
- Figure2.12:** Absorbance of two selected wavelengths plotted against each other page 20)
- Figure2.13:** Warping between two time series (page 25)
- Figure2.14:** Possible assignment between the vector pairs of  $\vec{x}$  and  $\vec{w}$  (page 28)
- Figure2.15:** A cost matrix with the minimum-distance warp path traced through it (page 29)
- Figure2.16:** The order that the cost matrix is filled (page 30)
- Figure2.17:** Flowchart of the online model of dynamic time warping algorithm (page 31)
- Figure2.18:** Two constraints: Sakoe-Chuba Band (left) and an Itakura Parallelogram (right), both have a width of 5 (page 33)
- Figure2.19:** Speeding up DTW by data abstraction (page 34)



**Figure2.20:** The four different resolutions evaluated during a complete run of the fast DTW algorithm (page 35)

### **List of Tables**

**Table3.1:** System performance regarding to gender and in silent room (page 37)

**Table3.2:** System performance regarding to gender in noisy room (page 37)

**Table3.3:** System performance regarding to age and gender in silent room (page 38)

**Table3.4:** System performance regarding to age and gender in noisy room (page 38)

**Table4.1:** Cost analysis of the project (page 38)