

# A Yet Efficient Communication System with Hearing-Impaired People Based on Isolated Words of Arabic Language

Khalid A. Darabkh, Ala F. Khalifeh, Iyad F. Jafar, Baraa A. Bathech, and Saed W. Sabah

**Abstract**— Although the mental capabilities of hearing-impaired people and hearing people are the same, reading and writing skills of those with hearing loss are much lower. Therefore, developing communications systems using sign languages becomes feasible and of interest. On the other hand, despite the fact that Arabic language is currently one of the most common languages worldwide, there has been only a little research on Arabic speech recognition relative to other languages such as English and Japanese. In this paper, a speech recognition system was developed which basically identifies the Arabic words that a hearing person speaks into a microphone to be then translated into a video, implements real Arabic sign language, which is easily understood by a hearing-impaired person. The speech recognition process carried out in this paper involves mainly using voice activity detection (VAD), Mel-frequency cepstral coefficients (MFCC), and dynamic time warping (DTW) algorithms. Moreover, delta and acceleration (delta-delta) coefficients have been added for the reason of improving the recognition accuracy. Utilizing the best set up made for all affected parameters to the aforementioned techniques, the proposed system achieved a recognition rate of about 98.5% which outperformed other relevant hidden Markov model (HMM) and artificial neural network (ANN)-based approaches available in the literature.

**Index Terms**— Arabic speech recognition, MFCC, DTW, VAD, delta and acceleration coefficients

## I. INTRODUCTION

THE phrase ‘hearing-impaired’ pertains to an individual with any degree of hearing loss or someone who cannot understand spoken words through hearing alone. Generally, the hearing loss may occur due to many different causes such as trauma, illness, genetics, old aging, and frequent exposing to loud voices [1]. The hearing loss may appear suddenly or

gradually. On the other hand, speaking loudly with hearing-impaired people rarely improves the conversation even when dealing with wearers of hearing aids since the distortion and pain are the results when amplifying the loud voice [2-3]. Hence, adopting efficient communication approaches are of interest and very helpful to interact with hearing-impaired people. It is noteworthy to mention that the communication with a hearing-impaired person can be performed in many forms, out of which facial expressions and gestures, speech reading, and sign language [2-5]. In the first form, associating the use of facial expression, gestures, as well as animation, which are affected by the mood, along with a spoken message at cultural bias intended times will be readily understood by hearing-impaired people. The second form, which is speech reading, includes lip reading to understand a spoken message and this can be done through watching lips’ movements. Furthermore, visual clues such as eye contact, pantomime, gestures, and rate of delivery can be incorporated to obtain a better enhancement. One drawback of speech reading form is the variation in hearing-impaired people’s ability to use speech reading which depends on each person’s vocabulary, amount of experience, practice in speech reading, and background noise interfering with his/her hearing [6-7]. Additionally, hearing-impaired people should concentrate intensively in order to obtain the advantage of this form [8-9]. However, in order to efficiently use this form, spoken words should be short and simple.

The last form of communication with a hearing-impaired person includes using sign language. This form is very popular and mainly depends on signs that are produced by the use of hand shapes, hand motions, hand locations, facial expressions, head motions, and body posture, which of course are all perceived visually [4, 10-11]. Examples of this form are American, French, Japanese, and Arabic sign languages. It is worth mentioning that the use of this form may not be efficient when communication evokes between hearing-impaired and hearing people who do not understand this form of language. Many solutions have been proposed for communicating efficiently with impaired people using sign language including the use of telecommunication devices [12], mobile text phones [13], e-learning projects [14-20], phonetics devices [21]. However, in the Middle East, there is a system called Arabic sign language translation system (ArSL-TS) which has been proposed to run on cell phones for translating Arabic typed text into Arabic sign language animations [22]. This system consists of four main phases described in the

Manuscript received January 16, 2013; revised July 10, 2013.

Khalid A. Darabkh is with the Department of Computer Engineering, The University of Jordan, Amman 11942, Jordan (phone: +962-77-9103900; e-mail: k.darabkeh@ju.edu.jo).

Ala F. Khalifeh is with the Department of Communication Engineering, German Jordan University, Amman 11180, Jordan (phone: +962 6 429 4112; email: ala.khalifeh@gju.edu.jo).

Iyad F. Jafar is with the Department of Computer Engineering, The University of Jordan, Amman 11942, Jordan (phone: +962-77-6743437; e-mail: iyad.jafar@ju.edu.jo).

Baraa A. Bathech and Saed W. Sabah are engineers graduated from the University of Jordan, Amman 11942, Jordan.

following order: cell phone player that is necessary for the user to insert the text and view the required animation, text translator tool that performs the text translation, matching tool that is performed between the result of text translator and sign animation database, and finally the output that represents the sign language which certainly is stored in the sign animation database.

It cannot be missed up that the Arabic language is considered nowadays as the fifth widely used language [23] as there are more than 200 million people speak this language. Unfortunately, the research efforts are still limited in comparison with other languages such as English and Japanese as far as the automatic speech recognition (ASR) is concerned. However, it is memorable that the Arabic digits (one-nine) are polysyllabic words while zero is a monosyllable word [24]. On the other hand, Arabic phonemes can be found of two categories, namely, pharyngeal and emphatic phonemes. These categories are only found in Semitic languages such as Hebrew [24-25]. However, automatic speech recognition has received a great deal of attention by many researchers for decades which basically allows a computer to recognize spoken words recorded by its microphone. Speech recognition is employed in a wide area of applications include interfacing with deaf people, home automation, healthcare, robotics, and much more. Actually, various approaches were adopted for speech recognition which are mainly found in three categories [26], Template-based such as dynamic time warping (DTW), neural network-based such as artificial neural networks (ANNs), and statistics-based such as hidden Markov models (HMMs).

Unlike other languages, Arabic language is characterized by having tremendous dialectal variety, diacritic text material, as well as morphological complexity which all in turn challenge the researchers in proposing highly accurate Arabic recognition system. In [27], a morphology-based language model was investigated for the use in a speech recognition system for conversational Arabic. In [28], the authors investigated the discrepancies between dialectal and formal Arabic in a speech recognition system utilizing morphology-based language model, automatic vowel restoration, as well as the integration of out of corpus language model. In [29], the authors reported the feasibility of using the automatic diacritizing Arabic text in acoustic model training for ASR. In [30], the authors attempted to use Carnegie Mellon University (CMU) Sphinx speech recognition system, which is one of the most robust speech recognizers in English, to develop an extension useful for Arabic language. However, more relevant research articles are discussed and compared with our work in the results and discussions section.

In this paper, we propose an efficient system applicable for Arab hearing-impaired community that allows only a direct communication between hearing and hearing-impaired persons in the shape of voice-to-video. The proposed system is very useful for both hearing-impaired people, who have below-average reading abilities for text, and hearing people who are illiterate and also not able to learn Arabic sign

language (ArSL). As far as the speech recognition is concerned, our proposed system is categorized as DTW-based speech recognition system which is applicable for isolated words of Arabic language. A brief summary of our system is as follows: A preprocessing is made for not only noise reduction, but also normalization. Moreover, speech/non-speech regions of the voice signal are detected using voice activity detection (VAD) algorithm. In addition, segmenting the detected speech regions into manageable and well-defined segments for the purpose of facilitating the upcoming tasks has been considered. As a matter of fact, the segmentation of speech can be practically divided into two types; the first one, which is employed in this paper, is called "*Lexical*", which divides a sentence into separate words, while the other type is called "*Phonetic*", which is based on dividing each word into phones. After the segmentation is excogitated, the Mel-frequency cepstral coefficients (MFCC) approach is adopted due to its robustness and effectiveness compared to other well-known feature extraction approaches like linear predictive coding (LPC) [31-32]. Moreover, delta and acceleration (delta-delta) coefficients have been added for the sake of improving the recognition accuracy. Finally, DTW is used as a pattern matching algorithm due to its speed and efficiency in detecting similar patterns [33-34]. Many experiments have been conducted to find the best parameters required to achieve the best efficient Arabic speech recognizer. Preliminary results of our proposed system, which exclude studying the impact of employing delta and acceleration (delta-delta) coefficients on the recognition accuracy, have been presented in [35].

The rest of the paper is divided into four further sections. Section II provides an overview of our system while Section III describes the proposed system broadly. Section IV presents our experimental results, observations, and discussions. Finally, Section V concludes our work whereas future directions to step this work are also provided.

## II. SYSTEM OVERVIEW

The overview of our system is shown in Fig. 1. The figure describes the process of translating the talk of a hearing person into a meaningful Arabic sign language video recognized by a hearing-impaired person. Communication starts by a hearing person saying a word which is recorded and stored in the system. This step includes preprocessing and segmenting the stored words. The spoken word is then analyzed whereas the information, belongs to it, is compared to a set of pre-analyzed database of Arabic words. Finally, the resulting word is used to pick the suitable Arabic sign language video, which corresponds to it, from a video database. The algorithms that illustrate hearing-to-hearing-impaired communication process are detailed shortly.

## III. THE PROPOSED SYSTEM

Hearing-to-hearing-impaired communication process consists of several stages as shown in Fig. 2. These stages are extensively detailed through the rest of sections as follows:

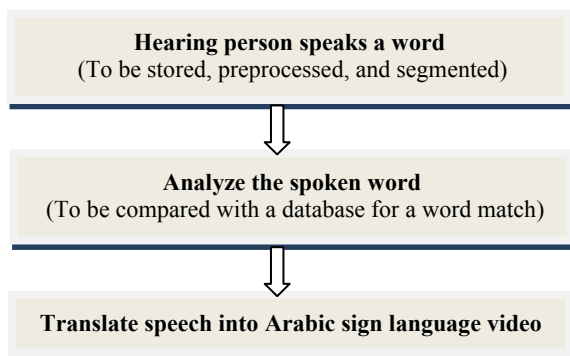


Fig. 1. Hearing-to-hearing-impaired communication system overview

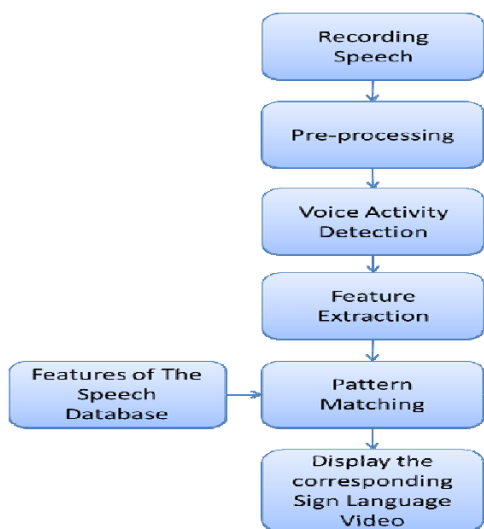


Fig. 2. Hearing-to-hearing-impaired communication stages

#### A. Database Collection

There is a need for a feature database that includes stored spoken words in Arabic for pattern matching process explained later. We have built a feature database that consists of a huge set of utterances (Arabic words and digits) for testing purposes produced by many different speakers (males and females) who were asked to record each word three times. An important point to mention is that words stored in our database were recorded in a normal home environment with a sampling rate of 8 KHz and 16 bit depth. This process is shown in Fig. 3 and illustrated well in *Algorithm 1*. Further details about what are mentioned in *Algorithm 1* can be found while discussing the subsequent subsections.

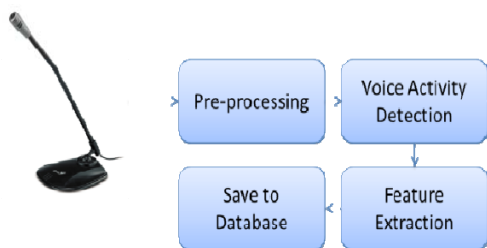


Fig. 3. Creating a feature database diagram

#### Algorithm 1: Pseudo code for creating the features database

```

//IN:
PATH = "Recorded Speech Path";
File;
//OUT:
DB = Database File;
Begin
    Generating Database;
    For File = 1:number of files in PATH do
        Read the File;
        Apply VAD;
        Extract Features;
        Save Features in DB;
    End for
End
    
```

#### B. Preprocessing

This stage aims to enhance some signal characteristics in order to achieve more accurate results through canceling disturbances that may affect the quality of the recorded speech. This stage is divided into two steps as follows:

##### Step 1: Pre-emphasis

At this step, high frequency contents of the input signal are emphasized in order to flatten the signal's spectrum. In our paper, the pre-emphasizer is represented by a first order FIR filter [36], which can be described according to:

$$H(z) = 1 - 0.97z^{-1} \quad (1)$$

where,  $z$  refers to discrete Fourier transform of the speech signal. However, the effect of applying this filter to a sample word in Arabic is shown in Fig. 4 whereas we can see that the high amplitude pulses in the signal that adversely affect the accuracy of stored word features in extraction stage are significantly reduced.

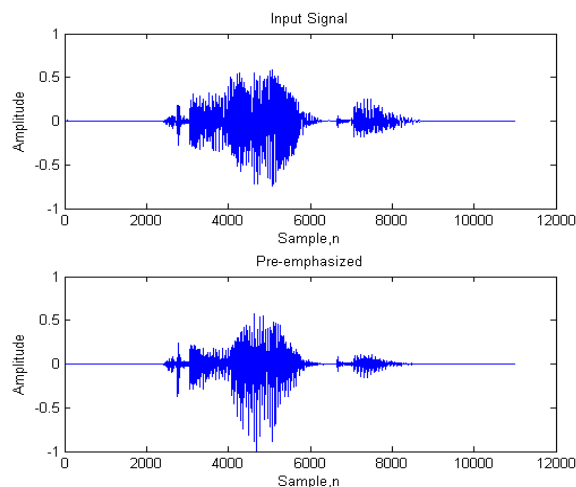


Fig. 4. The word "عليكم" after the pre-emphasis

##### Step 2: Hearingization

Speakers usually defer in speaking loudness [37]. Additionally, different microphones defer in their sensitivity

to speech. Thus, hearingization is included in our experiments which can be found as [38-39]:

$$S_1(n) = \frac{x_{pre-emphasis}(n) - \text{Mean}(x_{pre-emphasis}(n))}{\text{Max}(|x(n) - \text{Mean}(x_{pre-emphasis}(n))|)} \quad (2)$$

The hearingized version of the signal is depicted in Fig. 5.

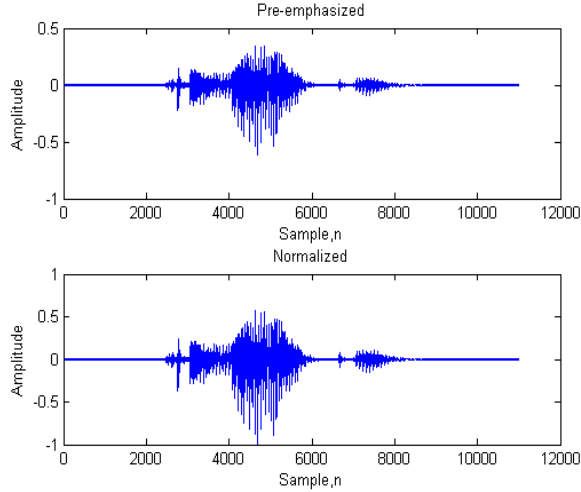


Fig. 5. The word “عليكم” after the hearingization

### C. Voice Activity Detection (VAD)

Generally, one of the major problems that affect the efficiency of a speech recognizer is detecting the start and end points of voice activity. However, short-term power and zero-crossing rate are commonly used parameters for distinguishing speech/non-speech regions [39]. Hence, this stage can be divided into the following steps:

#### Step 1: Framing

The speech signal is segmented into non-overlapped frames where each has a width of 20ms. Non-overlapping frames are used to reduce the number of times needed to check for voice activity. Consequently, the overall processing time of this stage is reduced.

#### Step 2: Short-term power and zero-crossing rate

It is worth mentioning that the short-term power is significantly increased in speech regions. However, it can be calculated according to [38]:

$$P_{S_1}(m) = \frac{1}{L} \sum_{n=m-L+1}^m S_1^2(n) \quad (3)$$

where  $m$ ,  $L$ , and  $n$  refer to frame number, frame length, frame index, respectively. On the other hand, zero-crossing rates tend to have larger values in non-speech regions. This gives a good indication of speech existence. Actually, it can be calculated using [39]:

$$Z_{S_1}(m) = \frac{1}{L} \sum_{n=m-L+1}^m \frac{|\text{sgn}(S_1(n)) - \text{sgn}(S_1(n-1))|}{2} \quad (4)$$

where,

$$\text{sgn}(S_1(n)) = \begin{cases} +1 & S_1(n) \geq 0 \\ -1 & S_1(n) < 0 \end{cases} \quad (5)$$

#### Step 3: Speech Indicator

The aforementioned parameters are combined in the following formula in order to provide a more comfortable approach which can be used to calculate a threshold value based on its mean and standard deviation [36, 39]:

$$W_{S_1}(m) = P_{S_1}(m)(1 - Z_{S_1}(m))F \quad (6)$$

where  $F$  is a constant scale factor used to avoid having small values. However, to initiate this function, we use the following activation function (AF):

$$AF_W = M_W + cSD_W \quad (7)$$

where,  $M_W$  refers to the mean of  $W_{S_1}(m)$ ,  $SD_W$  represents the standard deviation of  $W_{S_1}(m)$ , and  $c$  is a constant which should be fine-tuned since it depends on the signal characteristics. Accordingly, the voice activity detection function can be found as:

$$VAD(m) = \begin{cases} 1, & W_{S_1}(m) \geq AF_W \\ 0, & \text{Otherwise} \end{cases} \quad (8)$$

The result of this segmentation process is shown in Fig. 6. The output signal after performing VAD is  $x_1(n)$  where it is simply  $s_1(n)$  when  $VAD(n)$  is one. The steps of implementing VAD are described in *Algorithm 2*.

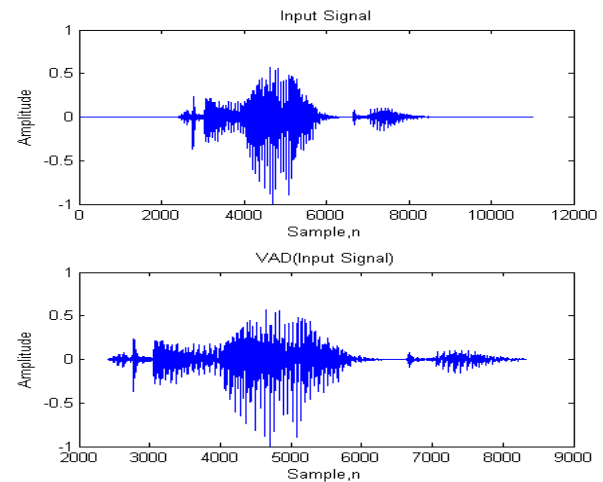


Fig. 6. The word “عليكم” as detected by VAD

#### Algorithm 2: Pseudo code for implementing VAD algorithm

```
//IN:
FrameLength = 160;
Overlap = 0%;
```

```

Signal;
//OUT:
VAD_Signal;
Begin
    Framed = Framing(Signal, FrameLength, Overlap);
    P = Power(Framed);
    Z = ZeroCrossing(Framed);
    W = P * (1-Z);
    M = Mean(W);
    S = Std(W);
    AF = M + c * S;
    VAD_Signal = (W>AF) * Signal;
End
    
```

#### D. Feature Extraction

The feature extraction phase consists of the following steps:

##### Step 1: Framing

In our experiments, the voice signal  $x_1(n)$  is broken up into  $J$  frames of  $P$  samples for each one with an overlapping ratio of 36.5% in which the adjacent frames are separated by  $T$  samples (where  $T < P$ ). The chosen values for  $P$  and  $T$  are 240 and 87 samples, respectively which were so appropriate. Hence, the output signal contains  $J$  vectors of length  $P$ , which corresponds to  $x_1(p; j)$ , where  $p$  and  $j$  vary from 0 to  $P-1$  and 0 to  $J-1$ , respectively.

##### Step 2: Hamming Window

Applying hamming window, to the output signal discussed in step 1 (framed signal), helps in reducing the discontinuity at both ends of each frame and this can be done utilizing the following formula [39]:

$$Ham(p) = 0.54 - 0.46 \cos \frac{2\pi p}{P-1}, \quad 0 \leq p \leq P-1 \quad (9)$$

where  $p$  refers to the sample index and  $P$  indicates the length of a frame (in samples). By applying  $Ham(p)$  to  $x_1(p; j)$  for all frames, then  $x_2(p; j)$ , which refers to the windowed signal, can be easily found.

##### Step 3: Fast Fourier Transform

To study the characteristics of the speech signal in frequency domain, we use  $N$ -point FFT to convert the windowed signal, resulting from step 2, from time domain to frequency domain. Note that the frame length here is a power of 2 ( $N=2^n$ ). Hence, the output signal is  $X_2(n; j)$ .

##### Step 4: Mel Filter Bank

According to the fact that human perception of voice frequencies is nonlinear (i.e., human hearing is less sensitive at higher frequencies, roughly  $> 1000$  Hz), a Mel-scale is used so that for each tone with a frequency  $F$  measured in Hz, a subjective pitch is measured on a Mel-scale according to following formula [36-38]:

$$F_{mel} = 2595 \log_{10} \left( 1 + \frac{F_{Hz}}{700} \right) \quad (10)$$

After finding the magnitude of  $X_2(n; j)$  and using the Mel scale filter bank (which consists of different triangular-band-pass filters that have an equal spacing before 1 KHz and logarithmic scale after 1 KHz), the Mel spectrum coefficients are found as the summation of the filtered results as the following:

$$Mel_v = \sum_{n=0}^{N-1} |X_2(n; j)| TF_v^{mel}(n) \quad (11)$$

where  $TF_v^{mel}(n)$  is the  $n^{\text{th}}$  triangular filter.

##### Step 5: Inverse Discrete Cosine Transform

To this end, we should return back to time domain. The best technique to do this while achieving highly uncorrelated features is the inverse discrete cosine transform (IDCT) as found in Equation (12). Before finding that, we compute first the logarithm of the magnitude of the output of Mel-filter bank since logarithm compresses dynamic range of values whereas humans are less sensitive to slight differences in the high amplitudes than low amplitudes.

$$IC(p; j) = \sum_{i=0}^{P-1} \lambda_i \log(Mel_i) \cos \left( \frac{\pi(2p+1)i}{2P} \right), \quad p = 0, 1, 2, \dots, P-1 \quad (12)$$

where,  $\lambda_0 = \sqrt{1/P}$  and  $\lambda_i = \sqrt{2/P}$ ,  $1 \leq i \leq P-1$

##### Step 6: Liftering

To extract the vocal tract cepstrum, it is good to use liftering which is mainly a filtering from frequency domain perspective. The simplest way to do that is to drop some of the cepstrum coefficient at the end. However, the most popular lifter that gives very promising recognition result is [39]:

$$l(p) = \begin{cases} 1 + \frac{Y-1}{2} \sin \left( \frac{\pi p}{Y-1} \right), & p = 0, 1, \dots, Y-1 \\ 0, & \text{Otherwise} \end{cases} \quad (13)$$

As a summary, we use the first 12 cepstral coefficients for each frame and ignore the rest which have F0 spike. In our work, the MFCC consists of steps 1 through 6.

##### Step 7: Short-term Energy

The cepstral coefficients do not capture energy. Therefore, the log of signal energy is an interesting feature to increase the coefficients derived from Mel-cepstrum. In other words, for every frame, the following energy term is added:

$$E_j = \log \sum_{p=0}^{P-1} x_2^2(p; j) \quad (14)$$

### Step 8: Delta and Acceleration Coefficients

It is known that the speech signal is not constant. In other words, the slope of formants may change from stop burst to release [40]. Hence, it is worth adding these changes in the features (i.e., the slopes). These are called delta features and delta acceleration (delta-delta) features. The delta coefficients are computed using a linear regression formula given  $2C+1$  is the size of the regression window:

$$\Delta IC_l(m) = \frac{\sum_{i=1}^C i(IC_l(m+i) - IC_l(m-i))}{2 \sum_{i=1}^C i^2} \quad (15)$$

where  $IC_l(m)$  is the  $m^{\text{th}}$  MFCC coefficient. As far as the delta-delta coefficients are concerned, they are found using linear regression of delta features. As a summary, we used 39-dimensional features as the following (12 MFCC, 1 energy feature, 12 delta MFCC features, 12 delta-delta MFCC features, 1 delta energy feature, 1 delta-delta energy feature). The steps of obtaining the features' vectors are described in Fig. 7. On the other hand, *Algorithm 3* describes the steps for analyzing the signal and obtaining its features along with delta and delta-delta coefficients.

---

#### **Algorithm 3: Pseudo code for extracting the features vector**

---

```

//IN:
VAD_Signal;
FrameLength = 240;
Overlap = 36.5%;
K = 20;
HammingWindow;
//OUT:
MFCCs;
Features_Extraction;
Begin
    Framed = Framing(VAD_Signal);
    Hammed = Filter(Framed, HammingWindow);
    FFT_Signal = FFT(Hammed);
    MelCep = MelFilterBank(ABS(FFT_Signal)^2, K);
    Cep = Log(IDCT(MelCep));
    Lifted = Liftering(Cep, 3:14);
    Energy = sum(Hammed^2);
    MFCC = [Energy:Cep];
    DeltaMFCC = Delta(MFCC);
    AccMFCC = Acceleration(MFCC);
    MFCCs = [MFCC:DeltaMFCC:AccMFCC]
End

```

---

### E. Pattern Matching

Dynamic time warping algorithm, which is based on dynamic programming, is a technique that calculates the level of similarity between two time series in which any of them

may be warped into a non-linear fashion by shrinking and stretching the time axis [32-34]. Fig. 8 shows the warp path between a tested word and stored spoken word as a result of applying DTW algorithm on different utterances. As shown from this figure, two time series are warped to find the best alignment between them. The lines shown between the two time series connect points that have the same value but happened in different time instants. Moreover, if the compared time series were identical, all lines connected between them must be straight vertical. Importantly, the warp path represents the actual distance between the two time series which can be measured as the accumulative sum between each two identical points in the time series being under comparison [33, 40]. To this extent, we can summarize that any tested word is segmented and its features are calculated and consequently compared with the whole database using DTW in order to find the word that has the nearest distance path to it. Thenceforth, a corresponding video, which describe the voice translation in Arabic sign language, will be displayed. This is clearly described in *Algorithm 4*.

## IV. RESULTS AND DISCUSSIONS

### A. Our results

In order to efficiently evaluate the performance of the proposed system, the minimum number of tests made to recognize an Arabic word is ten. Below is the formula which describes how the recognition rate of each word was calculated:

$$RR = \frac{\text{Number of correctly recognized words}}{\text{Number of tested words}} * 100\% \quad (16)$$

Table 1 describes the recognition rate for a sample of tested words in the database which we have previously recorded. For every tested word, the recognition rates are calculated using three different combinations of features as shown in this table. The positive effect of employing VAD and MFCC on the recognition rate is clearly observed. Furthermore, adding delta and acceleration coefficients to the feature set improves the recognition rate significantly. This is to be expected since the delta coefficients find the first derivative of the feature set which adds an important parameter that reflects the changing of speech from a specific phoneme to the next one. It is noteworthy to mention that the first derivative may give noisy results. Thus, in our proposed system, it is combined with the polynomial approximation approach. Consequently, the system's response for some tested words like "واحد" is improved. Additionally, by incorporating the polynomial approximation approach, it would be possible to calculate the second derivative of the features in order to provide more important information to the feature set for the reason of improving the overall performance and accuracy of the system. Actually, this can be noticed when considering the word "كم" shown in table 1. To this end, the fitting width ( $2C+1$ ) adds a delay of  $C$  blocks to the system. The choice of the value of  $C$  is a tradeoff between good and accurate approximation and long delay. Hence,  $C$  is given a value of 3 in order to have a good accuracy with a relatively faster response.



### B. Comparisons with previous work

There are interesting approaches, similar in target to our proposed system, done to improve the recognition rate of Arabic language. As mentioned previously concerning [30], an Arabic speech recognition system was proposed using open source CMU Sphinx-4 and hidden Markov models. The obtained recognition accuracy was about 85.55%. While in [32], a comparison of DTW and discrete hidden Markov model was made for recognizing isolated words in Arabic language. In DTW-based approach, they used 13 MFCC coefficients and also the same for delta and acceleration coefficients. A 256 point FFT was used to find the power spectrum to be used in an emulated filter-bank composed of 24 triangular weighting functions in Mel scale. Thereafter, the natural logarithmic was applied to the 24 filter-bank. They measured the recognition rate for frames' overlapping length of 512\*256. The recognition accuracy was about 86% in clear environment using the characteristics of power (energy) and differential information ( $\Delta$  and  $\Delta\Delta$ ). In DHMM-based speech recognizer, five states were defined for each word whereas transitions between these states were possible only in left to right direction with no states skipping. No more details were reported about these states and transitions. The achieved recognition accuracy was about 92%.

In [41], a heuristic method for Arabic speech recognition (ArSR), minimal eigenvalues algorithm, was used to find the most promising path through a tree of different samples of an uttered word. Furthermore, radial neural networks (RNN) approach was incorporated with this heuristic method to enhance the recognition rate. The recognition accuracies were about 86.45% and 95.82% for minimal eigenvalues algorithm and RNN, respectively. In [42], comparisons were made between monophone, triphone, syllable, and word-based algorithms for recognizing Egyptian Arabic digits. Thirty-nine MFCC coefficients were extracted as features for every recorded voice in the database where they were used to train HMMs in which the system matched between the testing word and training database. The achieved recognition accuracies were about 90.75%, 92.24%, 93.43%, and 91.64% for monophone, triphone, syllable, and word-based recognition algorithms, respectively. In [43], an Arabic numeral recognition (ArNR) technique was proposed using vector quantization (VQ) and HMM whereas the LP cepstral coefficients were used. The recognition accuracy was about 91%. In [44], HMM-based Arabic numeral recognition system was proposed using Wavelet cepstral coefficients and Mel frequency cepstral coefficients whereas the recognition accuracies for different numerals were about 61%-92% and 76%-92% for MFCC-based and Wavelet-based systems, respectively. In [45], other HMM-based approaches were proposed, based on LPC and MFCC, for Arabic numeral recognition and devised for field programmable gate arrays (FPGAs) whereas the recognition accuracy ranged from 91% to 96% for LBC-based recognition systems and 95% to 98% for MFCC-based recognition systems. The recognition rates obtained from aforementioned approaches are summarized in table 2. All mentioned rates are obtained assuming clear

environment. The significance of our proposed work is consistently noticed.

### V. CONCLUSIONS AND FUTURE WORK

Many solutions were proposed using sign languages to deal efficiently with hearing-impaired people which were focused on English or Japanese sign languages. However, finding powerful solutions that include Arabic sign language still lacks the research importance it deserves. In this paper, we have proposed a system that eases the communication between hearing people and people with hearing difficulties (in the direct way only) based on voice-to-video style. On the other hand, finding efficient automatic speech recognition techniques for Arabic words is of great interest since the research efforts remain limited. In this work, the robustness of MFCC combined with DTW algorithm is very conspicuous. Moreover, the voice activity detector technique has a significant impact on the system's performance. On the other hand, adding delta and delta-delta coefficients help much in improving the overall recognition accuracy. Many experiments were conducted to choose the best parameters that maximize the improvements of Arabic speech recognition. Additionally, a noticeable speech recognition accuracy improvement is achieved when compared to other HMM and ANN-based approaches. There are many future directions to this work. It is of interest to work on the reverse direction of communication in which the system can automatically convert Arabic sign language videos back into speech to be recognized by hearing people. Many interesting approaches can be adopted to achieve this goal, out of which, recognition of gestures using neuro-fuzzy systems, boosted hidden Markov models, object tracking using snake algorithms, and artificial neural networks combined with graph matching algorithm.

### REFERENCES

- [1] Warren R. Goldmann, James R. Mallory "Overcoming Communication Barriers: Communicating with Deaf people," *Library Trends: University of Illinois*, vol. 41, no.1, pp. 21-30, Summer 1992.
- [2] Susan V. Rezen, Carl Hausman, *Coping With Hearing Loss: A Guide for Adults and Their Families*, Barricade Books Inc., 1993.
- [3] Donald Moores, *Educating the Deaf: Psychology, Principles, and Practices*, Wadsworth Publishing, August 2000.
- [4] Alice Caplier, Sébastien Stillitano, Oya Aran, Lale Akarun, Gérard Bailly, Denis Beutemps, Nouredine Aboutabit, and Thomas Burger, "Image and Video for Hearing Impaired People," *EURASIP Journal on Image and Video Processing*, Volume 2007, Article ID 45641, 14 pages, December 2007.
- [5] David Lodge, *Deaf Sentence*, Penguin Books, June 2009.
- [6] N. Eveno, A. Caplier, and P.-Y. Coulon, "Automatic and accurate lip tracking," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 14, no. 5, pp. 706-715, 2004.
- [7] N. Aboutabit, D. Beutemps, and L. Besacier, "Automatic identification of vowels in the Cued Speech context," in *Proceedings of the International Conference on Auditory-Visual Speech Processing (AVSP '07)*, Hilvarenbeek, The Netherlands, August 2007.

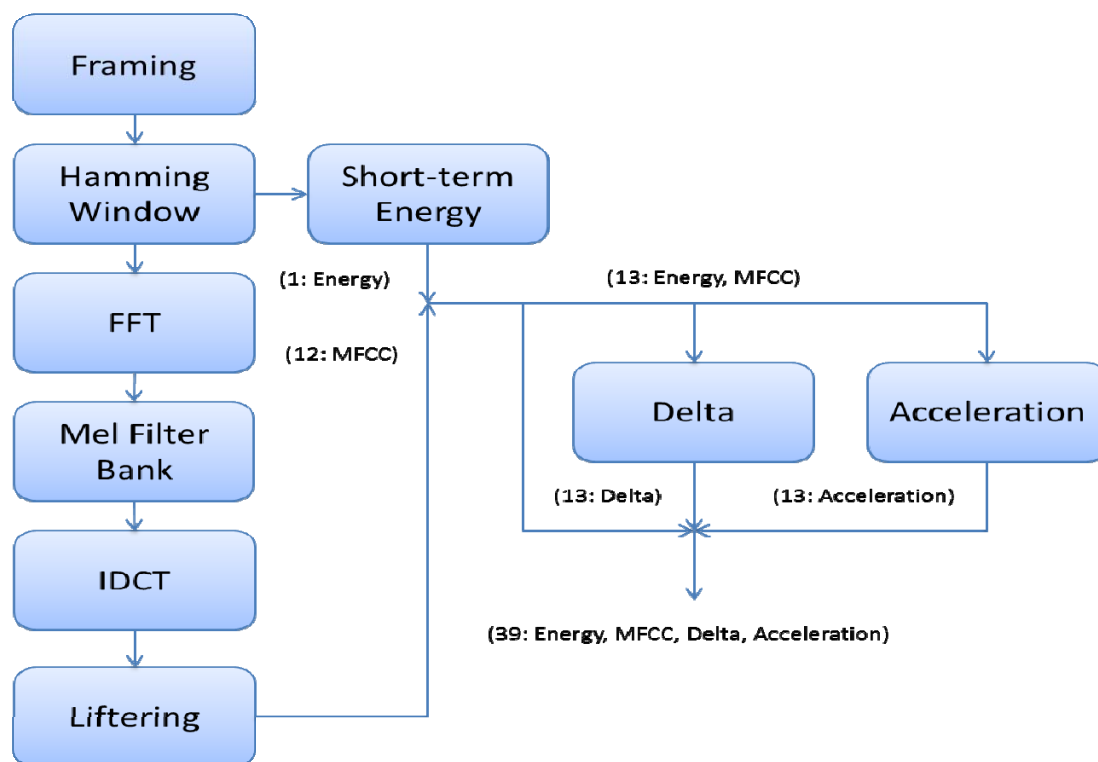


Fig. 7. Features extraction processing steps

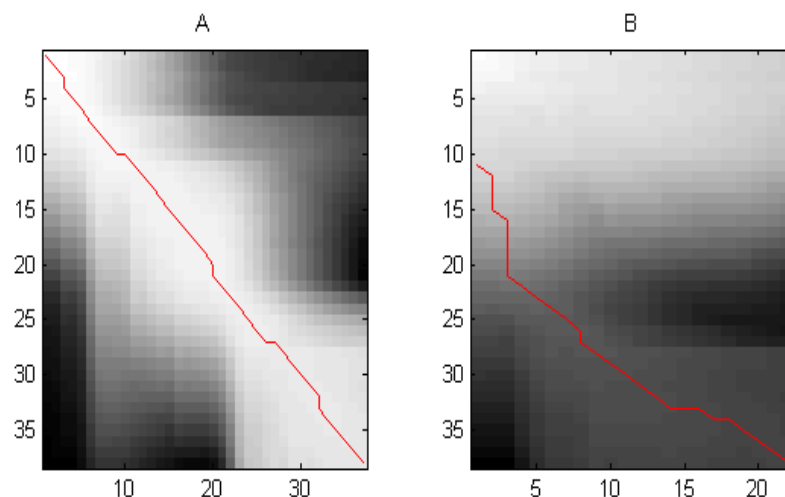


Fig. 8. The warp path after applying DTW algorithm: A) using different utterances of the word "عليكم", B) using the words "عليكم" and "سبعة"

**Algorithm 4: Pseudo code for pattern matching**

```

//IN:
RecordedSpeech;
DB;
//OUT:
Displayed_Video;
Begin
VAD_Speech = VAD(RecordedSpeech);
Features = FeatureExtraction(VAD_Speech);
Cost = Infinity;
Name;

```

```

Matching;
For N = 1:Length(DB)
    SM = SimilarityMatrix(Features, Feature(DB(N)));
    CurrentCost = dpfast(1-SM);
    If CurrentCost < Cost
        Cost = CurrentCost;
        Name = Name(DB(N));
    End if
End for
DisplayVideo(Name);
End

```



**TABLE I.**  
RECOGNITION RATES FOR DIFFERENT FEATURE SETS

RECOGNITION RATES					
Tested Word (Arabic Writing)	Transcription	English Writing	Approach 1: VAD+MFCC	Approach 2: VAD+MFCC+Δ	Approach 3: VAD+MFCC+Δ+ΔΔ
واحد	WAHID	ONE	85.7%	100%	100%
اثنان	ITHNAN	TWO	100%	100%	100%
ثلاثة	THALATHA	THREE	100%	100%	100%
أربعة	ARBAA	FOUR	100%	100%	100%
خمسة	KHAMSA	FIVE	100%	100%	100%
ستة	SITTA	SIX	85.7%	85.7%	85.7%
سبعة	SABAA	SEVEN	100%	100%	100%
ثمانية	THAMANIYA	EIGHT	100%	100%	100%
تسعة	TISAA	NINE	100%	100%	100%
عشرة	ASHRA	TEN	85.7%	100%	100%
السلام	ASSALAAMU	PEACE	100%	100%	100%
عليكم	ALAIKUM	UPON YOU	100%	100%	100%
كيف	KEEF	HOW	100%	100%	100%
حالك	HALAK	ARE YOU	85.7%	85.7%	85.7%
ما	MA	WHAT	100%	100%	100%
اسمك	ESMOK	YOUR NAME	100%	100%	100%
كم	KAM	HOW	85.7%	85.7%	100%
عمرك	OMROK	YOUR AGE	100%	100%	100%
مهنتك	MEHNATOK	YOUR OCCUPATION	100%	100%	100%

**TABLE II.**  
COMPARISONS WITH PREVIOUS WORK

Previous Work	Recognition Rates
ASR using CMUSphinx [30]	85.55%
DTW-Based ArSR [32]	86%
DHMM-Based ArSR [32]	92%
Heuristic Method [41]	86.45%
Heuristic Method with RNN [41]	95.82%
Monophone-Based ArSR [42]	90.75%
Triphone-Based ArSR [42]	92.24%
Syllable-Based ArSR [42]	93.43%
Word-Based ArSR [42]	91.64%
VQ and HMM ArNR [43]	91%
MFCC-based ArNR [44]	61%-92%
Wavelet-based ArNR [44]	76%-92%
LBC-based FPGA ArNR [45]	91%-96%
MFCC-based FPGA ArNR [45]	95%-98%
Our proposed recognition system	98.5%

- [8] V. Attina, D. Beautemps, M.-A. Cathiard, and M. Odisio, "A pilot study of temporal organization in Cued Speech production of French syllables: rules for a Cued Speech synthesizer," *Speech Communication*, vol. 44, no. 1–4, pp. 197–214, 2004.
- [9] S. Stillitano and A. Caplier, "Inner lip segmentation by combining active contours and parametric models," in *Proceedings of the 3rd International Conference on Computer Vision Theory and Applications (VISAPP '08)*, Madeira, Portugal, January 2008.
- [10] S. K. Liddell, *Grammar, Gesture, and Meaning in American Sign Language*, Cambridge University Press, Cambridge, UK, 2003.
- [11] W. C. Stokoe Jr., "Sign language structure: an outline of the visual communication systems of the American deaf," *Studies in Linguistics: Occasional papers* 8, Buffalo, NY: Department of Anthropology and Linguistics, University of Buffalo, 1960 (Reprinted in *Journal of Deaf Studies and Deaf Education*, vol. 10, no. 1, pp. 3–37, 2005).
- [12] Virginia W. Stern and Martha Ross Redden, "Selected Telecommunications Devices for Hearing-Impaired Persons," *Office of Technology Assessment*, <http://www.fas.org/ota/reports/8225.pdf>.
- [13] Sean Anatasi, Michael Eng, Joy Kim, Rafael Rodriguez, Sherri Yin, *MobileASL*, University of Washington, <http://mobileasl.cs.washington.edu>.
- [14] Michael S. Stinson, Lisa B. Elliot, Ronald R. Kelly, and Yufang Liu, "Deaf and hard-of-hearing students' memory of lectures with speech-to-text and interpreting/note taking services," *The Journal of Special Education*, Hammill Institute on Disabilities, vol. 43, no. 1, pp. 52–64, May 2009.

- [15] Eiichi Takada (Japanese Federation of the Deaf), "Solidarity and Movements of the Deaf and Hard of Hearing in Asia," *Disability Research Magazine*, vol. 28, issue 1, 2000.
- [16] Giuliano Pirelli (European Commission Joint Research Centre), "The Voice Project: Giving a Voice to the Deaf by Developing Awareness of Voice to Text Recognition Capabilities," *Proceedings of the 1998 TIDE Conference*, Helsinki, Finland, July 1998.
- [17] Katja Straetz, Andreas Kaibel, Vivian Raithel, Marcus Specht, Klaudia Grote, and Florian Kramer, "An e-Learning Environment for Deaf Adults", *Proceedings of the 8th ERCIM Workshop on "User Interfaces for All"*, Vienna, Austria, June 2004.
- [18] A.S. Drigas and D. Kouremenos, "An e-Learning Management System for the Deaf people," *WSEAS Transactions on Advances in Engineering Education*, vol. 2, issue 1, pp. 20-24, 2005.
- [19] Meryl Glaser and William Tucker, "Telecommunications bridging between Deaf and hearing users in South Africa," *Proceeding of Conference and Workshop on Assistive Technologies for Vision and Hearing Impairment (CVHI 2004)*, Granada, Spain, June 2004.
- [20] A. S. Drigas, J. Vrettaros, and D. Kouremenos, "E-learning Environment for Deaf people in the E-Commerce and New Technologies Sector," *WSEAS Transactions on Information Science and Applications*, vol. 1, issue 5, November 2004.
- [21] Venkatraman.S, and T.V. Padmavathi, "Speech for the Disabled," *Proceedings of the International MultiConference of Engineers and Computer Scientists 2009 (IMECS 2009)*, vol. I, Hong Kong, March 2009.
- [22] Sami M. Halawani, "Arabic Sign Language Translation System on Mobile Devices," *International Journal of Computer Science and Network Security (IJCSNS)*, vol.8, no.1, 6 pages, January 2008.
- [23] M. Al-Zabibi, "An Acoustic-Phonetic Approach in Automatic Arabic Speech Recognition," *The British Library in Association with UMI*, UK, 1990, <http://hdl.handle.net/2134/6949>.
- [24] M. Alkhoul, "Alaswaat Alaghawaiyah," *Daar Alfalah*, Jordan, 1990 (in Arabic).
- [25] M. Elshafei, "Toward an Arabic Text-to-Speech System," *The Arabian Journal for Science and Engineering*, vol. 16, no. 4B, pp. 565-83, October 1991.
- [26] Patricia Melin, Jerica Urias, Daniel Solano, Miguel Soto, Miguel Lopez, and Oscar Castillo, "Voice Recognition with Neural Networks, Type-2 Fuzzy Logic and Genetic Algorithms," *Engineering Letters*, volume 13, no. 2, August 2006.
- [27] D. Vergyri, K. Kirchhoff, K. Duh, A. Stolcke, "Morphology-based language modeling for Arabic speech recognition", *In INTERSPEECH-2004*, pp. 2245-2248, 2004.
- [28] K. Kirchho, J. Bilmes, J. Henderson, R. Schwartz, M. Noamany, P. Schone, G. Ji, S. Das, M. Egan, F. He, D. Vergyri, D. Liu, and N. Duta, "Novel Approaches to Arabic Speech Recognition," *Technical Report*, Johns-Hopkins University, 2002.
- [29] D. Vergyri, K. Kirchhoff, "Automatic diacritization of Arabic for acoustic modeling in speech recognition", In Ali Farghaly and Karine Megerdumian, editors, *COLING 2004, Computational Approaches to Arabic Scriptbased Languages*, pp. 66-73, Geneva, Switzerland, 2004.
- [30] H. Satori, M. Harti, N. Chenfour, "Introduction to Arabic Speech Recognition Using CMUSphinx System," *Proceedings of Information and Communication Technologies International Symposium (ICTIS'07)*, Fes, Morocco, pp. 139-115, July 2007.
- [31] S.B. Davis, P. Mermelstein, "Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 28, no.4, pp. 357-366, August 1980.
- [32] Z. Hachkar, A. Farchi, B. Mounir, J. El Abbadi, "A Comparison of DHMM and DTW for Isolated Digits Recognition System of Arabic Language," *International Journal on Computer Science and Engineering*, vol.3, no.3, pp.1002-1008, March 2011.
- [33] Lindasalwa Muda, Mumtaj Begam, I. Elamvazuthi, "Voice Recognition Algorithms using Mel Frequency Cepstral Coefficient (MFCC) and Dynamic Time Warping (DTW)", *Journal of Computing*, vol. 2, no. 3, pp. 138-143, March 2010.
- [34] Stan Salvador, Philip Chan, "Toward Accurate Dynamic Time Warping in Linear Time and Space", *Intelligent Data Analysis Journal*, vol. 11, no. 5, pp. 561-580, October 2007.
- [35] K. Darabkh, A. Khalifeh, I. Jafar, B. Bathech, and S. Sabah, "Efficient DTW-Based Speech Recognition System for Isolated Words of Arabic Language," *Proceedings of International Conference on Electrical and Computer Systems Engineering (ICECSE 2013)*, Lucerne, Switzerland, pp. 699-702, May 2013.
- [36] Lawrence Rabiner, Biing-Hwang Juang, *Fundamentals of speech recognition*, Upper Saddle River, New Jersey: Prentice Hall, USA, 1993
- [37] X. Huang, A. Acero, and H.-W. Hon, *Spoken Language Processing*, Upper Saddle River, New Jersey: Prentice Hall, USA, 2001.
- [38] B. Gold and N. Morgan, *Speech and Audio Signal Processing*, New York, New York: John Wiley and Sons, USA, 2000.
- [39] Mikael Nilsson and Marcus Einarsson, "Speech Recognition using Hidden Markov Model (performance evaluation in noisy environment)", *Masters Thesis*, Department of Telecommunications and Signal Processing, Belkinge Institute of Technology, Ronneby, Sweden, March 2002.
- [40] B.S. Jinjin Ye, "Speech Recognition Using Time Domain Features From Phase Space Reconstructions", *Masters Thesis*, Department of Electrical and Computer Engineering, Marquette University, Milwaukee, Wisconsin, May 2004.
- [41] Khalid Saeed and Mohammad Nammous, Heuristic Method of Arabic Speech Recognition, *Bialystok University of Technology*, Poland, <http://aragorn.pb.bialystok.pl/~zspinfo/>
- [42] Mohamed Mostafa Azmi, Hesham Tolba, Sherif Mahdy, Mervat Fashal, "Syllable-Based Automatic Arabic Speech Recognition", *Proceedings of WSEAS International conference of Signal Processing, Robotics and Automation (ISPRA' 08)*, University of Cambridge, UK, pp. 246-250, February 2008.
- [43] H. Bahi and M. Sellami, "Combination of Vector Quantization and Hidden Markov Models for Arabic Speech Recognition," *Proceedings of the ACS/IEEE International Conference on Computer Systems and Applications (AICCSA 2001)*, Beirut, Lebanon, pp: 96-100, June 2001.
- [44] W. Alkhaldi, W. Fakhr, N. Hamdy, "Multi-Band Based Recognition of Spoken Arabic Numerals Using Wavelet Transform," *Proceedings of the 19<sup>th</sup> National Radio Science Conference (NRSC'01)*, Alexandria University, Alexandria, Egypt, March 19-21, 2002.
- [45] F.A. Elmisery, A.H. Khalil, A.E. Salama, H.F. Hammed, "A FPGA Based HMM for a Discrete Arabic Speech Recognition System," *Proceedings of the 15<sup>th</sup> International Conference on Microelectronics (ICM 2003)*, Cairo, Egypt, December 9-11, 2003.