# Sound Visualization for the Hearing Impaired

Jimmy Azar, Hassan Abou Saleh, and Mohamad Adnan Al-Alaoui

American University of Beirut/ECE Department, Beirut, Lebanon

*Abstract*—In this paper, we investigate several means of visualizing both ambient and speech sounds and present a fusion of different visualization displays into one program package that would help provide the hearing impaired with a means to an enhanced awareness of their surroundings. The ideas investigated were implemented in software, and the program was evaluated by means of a survey conducted in a school for the deaf.

*Index Terms*—Ambient Sounds, Autocorrelation Function, Center Clipping, Pitch, Sound Visualization, Spectrogram, Speech Recognition.

## I. INTRODUCTION

There are approximately 219,480 (5.487%) deaf people in Lebanon [1] and nearly 10,000,000 persons with hearing impairments and close to 1,000,000 who are functionally deaf in the United States [2]. Without Assistive Technologies, there is no possibility for the hearing impaired to recognize sounds efficiently.

Medical or surgical solutions such as cochlear implants may not always be possible. Methods such as mapping frequency and intensity of sound to the frequency and intensity of vibrations (as in an alarm pillow or vibrating movie seat) are limited compared to what visual displays may offer. Other means involving flashing lights remain limited due to their dependence upon some prior knowledge of a characteristic of the specific sound signal expected. The Positional Ripples Display [3] as a proposed method of conveying information of amplitude and location has its limitations in that the physical setup (array of distributed microphones, etc.) may be involved and remains bound to a stationary workplace requiring prior knowledge of the architectural setup. The idea of spectrogram visualization [3-4] is further emphasized in this paper to reveal the capability of such a display in pattern visualization. Furthermore, other sound to image mappings were implemented in addition to speech recognition and Speech-to-ASL (American Sign Language) translation. By this, we would have provided a comprehensive means of acquiring information from verbal and non-verbal sounds.

### A. Design Overview

Several sub-windows were designed with each providing specific sound information. All these sub-windows are dynamically changing as sound is continuously inputted (live), and they are all simultaneously displayed within a Main Window. The sub-windows may be divided into those visualizing information of ambient sounds and those displaying speech. Belonging to the first category are the Dynamic Circles Display (for visualizing pitch, sound location, and loudness), Dynamic Spectrogram Display (for pattern visualization), Pitch Vector Display (for speech visualization and emphasis of pitch presence/absence), RMS Line Display, and RMS Image Display (for visualizing changes in sound intensity via a plot or color strip image). Belonging to the second category are the Speech Recognition Display and the Speech-to-ASL Display.

### B. System Setup

The program was written using MATLAB version 6.1 and run on a 1.7 GHz Intel Pentium processor. The data acquisition system consisted simply of a pair of identical microphones introduced into the 'Line in' socket. The system was adjusted for 'Line in' Recording. Sound signals were inputted via two separate channels each corresponding to a microphone. Hence, the input consists of two separate sound signals which may be processed separately in the case of FFT Left/Right displays as well as in the Dynamic Circles display, but are combined to present displays for the other sub-windows.

Section II of the paper discusses the various means of visualizing ambient sounds. Different displays are proposed and implemented.

Section III discusses speech visualization which includes the implementation of speech recognition and speech to ASL translation followed by an overall view of the entire Main Window Display.

Section IV includes a summary of the statistical results gathered from a survey conducted in a school for the hearing impaired with the intention of assessing the usefulness and performance of our program. Possible applications are also discussed.

Section V contains the conclusion of our work and suggestions for possible future improvements.

## II. AMBIENT SOUND VISUALIZATION

### A. Pitch Extraction Algorithm

The Autocorrelation Function (ACF) as defined in (1)

$$R_y[k] = < y[n]\, y[n-k] > \tag{1}$$

(where <.> corresponds to the mean) is essential for extracting the pitch value of a sound segment. Prior to extracting the pitch, preconditioning is applied mainly in the form of Center Clipping and median filtering. Pitch contributes to significant peaks in the Autocorrelation function (ACF), however in between there may exist secondary peaks corresponding to the vocal tract transfer function. These peaks may cause errors in the automated pitch detection process. Center Clipping the speech signal by considering only the peaks that overpass 68% of the minimum or maximum values of the signal eliminates these undesirable peaks and hence allows for a more accurate detection of pitch [5].
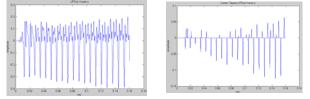


Figure 1: (a) LPF'ed Vowel Segment('a');   (b) Effect of Center Clipping

Median Filtering is then applied over the speech segment before computing the ACF in order to remove outlier components outside the human voice range. The sound segment is then divided into intervals of 50 milliseconds overlapping midway. To avoid unnecessary peaks, the segments are further Low Pass filtered using a finite impulse response (FIR) filter of order 125 with a cutoff frequency equal to 500 Hz since the upper limit of the human pitch range is ~320 Hz. Finally, the ACF is processed for each sub-segment to yield a pitch value for each sub-segment.
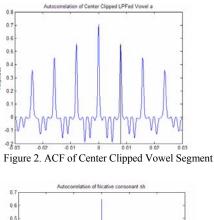
Pitch may be extracted from the ACF by means of a simple algorithm described below:

(i) The positive part of the symmetric autocorrelation function (i.e. from Ry(0) →end) is considered.

(ii) Identify/Find the index at which the first peak (that of Ry(0) ) actually ends.

(iii) Search for the 2$^{nd}$ peak 'pk' which is just the maximum value of the remaining part of the function starting from the discovered index in (ii).

(iv) The correct index of the second peak is then converted to real time by multiplying with the sampling period Ts to obtain 'pk_t'.

(v) To be considered a valid peak, the second peak 'pk' has to be greater than 30% of the first peak Ry(0).

The pitch is then :

$$\begin{cases} \dfrac{1}{pk\_t}, & \text{if } pk > 0.3\, Ry(0) \\[2mm] 0, & \text{Otherwise} \end{cases}$$

For instance pitch is equal to 0 for consonant 'sh' which is of an unvoiced fricative nature as compared to a stressed vowel segment like 'a'.
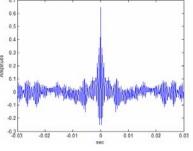


Figure 2. ACF of Center Clipped Vowel Segment



Figure 3. ACF of consonant 'sh'

### B. Pitch Vector Display

The presence of definitive pitch is often indicative of intelligible sounds. For example, pitch is apparent in speech mainly due to vowels and sound originating from the vocal cord (e.g. 'zzz') as well as in musical instruments which possess a large range of pitch values. Fricative or consonant unvoiced sounds which in the case of speech usually result from teeth, tongue, and lips (e.g. 'sss') do not possess a pitch value. An incoming sound segment is further divided into 50 msec intervals overlapping midway, and the 'detect_pitch' function is used to obtain the corresponding pitch value for each of these sub-segments. The Pitch Vector Display is basically a plot of these values after normalization. The normalized pitch values are plotted (in black) on top of the speech signal segment (in blue) in a sub-window in between the Dynamic Circles and Dynamic Spectrogram sub-windows. The idea is to display the presence of pitch and its duration, but more importantly to give a clearer view of the black pitch line as indicative of pitch activity. This sub-window is to be considered complimentary to the two other displays in that it further emphasizes the presence or absence of pitch. The Pitch Vector is displayed coincided with a plot of the acquired normalized speech segment with the dc mean removed (Fig. 4).
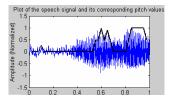
Figure 4. Pitch Vector Display (black line) with normalized sound segment overlapping

## C. Dynamic Circles Display

When asked to depict sound, ripples or circles are usually adopted [3]. Hence, we developed The Dynamic Circles Display which combines sound localization and pitch extraction. The location of the center of the circle is indicative of the location of sound linearly between the two microphones and the mean pitch extracted from the acquired sound segment. The physical location between left and right is mapped to the x-location of the center, and the mean pitch value is mapped to the y-location. The pitch vector discussed often contains several 0 elements which are then removed and the remaining values averaged to yield a mean pitch value for the entire segment—this is the y-ordinate of the center of the circle. By this, non-pitch segments are ignored and the mean pitch corresponds to that of the pitch-possessive segments. If the vector is entirely a zero vector, then the mean pitch value is 0. The loudness of the sound segment is mapped simultaneously to the radius, color, and thickness of the circle. A new circle is displayed every second with every incoming sound segment. The window is cleared every 5 seconds so as to display a temporary history of 5 circles. In the case of speech, the range in which the value falls is indicative of the gender and/or age of the speaker as males have pitch values in the range [80 160] Hz while females and children in the range [160 320] Hz. The display of the loudness of sound may be indicative of the emotional state of the speaker (whether angry, calm, etc).
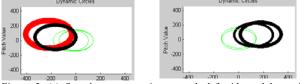


Figure 5. (a) Sound source moving towards left side and becoming louder; (b) Sound source moving towards right side and becoming louder. Pitch values for both is in the range [0 200]Hz.
Color Range: Green(low)→Blue→Black→Red(high)

## D. Dynamic Spectrogram

A spectrogram is a very useful visual tool for displaying the intensity of frequency components over time segments via color. With sound being highly non-stationary, the Short Time Fourier Transform (STFT) was used in this regard. This display reveals dominant frequency components (mainly formant band striations due to resonance of the vocal tract), their harmonics, and a sense of how the intensity of such components is changing with time. If well trained, words may be read from spectrograms based on the patterns that appear. This however is a difficult task and may be compared to learning a second language. The intention behind the spectrogram display is pattern visualization of events that

tend to recur. Examples of such events are phone ringing, door knocking, and other ambient or peripheral sounds.
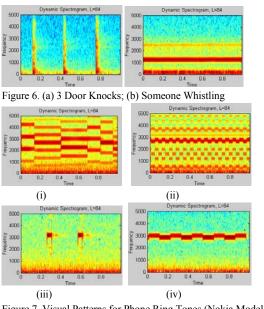


Figure 6. (a) 3 Door Knocks; (b) Someone Whistling



(i)            (ii)



(iii)           (iv)

Figure 7. Visual Patterns for Phone Ring Tones (Nokia Model 8850)
(i) 'Nokia' Tone;         (ii) Traditional 'low' tone;
(iii) 'Knock-Knock' Tone;     (iv) 'Mosquito' tone

The spectrogram used is wide-band with a frequency resolution of 300 Hz. The window utilized within the spectrogram is a Hanning window [6] described by (2)

$$w(n) = \frac{1}{2}\left[1 + \cos\left(\frac{2\,x\,\pi\,x\,n}{2M+1}\right)\right], \quad -M \le n \le M \tag{2}$$

It should be noted that the tradeoff between time and frequency resolution is described by the uncertainty principle [7] mainly stating:

If an impulse response h(t) treated as a probability density function (PDF) decays to zero faster than $\frac{1}{\sqrt{|t|}}$ for large t, then

$$(\Delta T)\,(\Delta\omega) \ge 2 \tag{3}$$

Where $(\Delta T)$, the time duration, is defined as twice the standard deviation $\sigma_h$ from its zero centered mean $\mu_h$:

$$(\Delta T)^2 = 4\,\sigma_h{}^2 = 4\frac{\int_{-\infty}^{\infty} t^2\,h(t)^2\,dt}{\int_{-\infty}^{\infty} h(t)^2\,dt} \tag{4}$$

and $(\Delta\omega)$, the bandwidth of $H(\omega)$, is defined as twice the standard deviation $\sigma_H$ of the PDF $|H(\omega)|^2$.

$$(\Delta\omega)^2 = 4\,\sigma_H{}^2 = 4\frac{\int_{-\infty}^{\infty} \omega^2\,|H(\omega)|^2\,d\omega}{\int_{-\infty}^{\infty} |H(\omega)|^2\,d\omega} \tag{5}$$

### E. Left / Right FFT Display

The input consists always of two sound segments acquired from the left and right microphones. The Discrete Fourier Transform (DFT) is applied on each of the two segments using the Fast Fourier Transform (FFT) algorithm. The magnitude of the DFT is then plotted versus frequency for the left and right signals separately. This display allows visually detecting the dominant and non-dominant frequency components present and comparing these between left and right sides.


(a)                                    (b)

Figure 8. Singing a high pitch 'aaa' from right side (a) DFT Display (left mic); (b) DFT Display (right mic) shows dominant frequency component in range [250 300] Hz.

### F. RMS Line Display

The root mean square value of a sound segment is a direct indication of loudness. The speech segment is divided into sub-intervals, and the root mean square value is calculated for each sub-interval. The result is then plotted to reveal how loudness is varying. This is of course complimentary to the thickness/color/and radius of the dynamic circle display with the difference that the rms display has a faster, clearer, and more accurate response to changes in loudness levels due to the division of sound into sub-intervals.



Figure 9. 'EEEE' sound showing exponentially rising loudness intensity.

### G. RMS Image Display

This display is an alternative to the RMS Line Display. Instead of plotting the rms values versus time, each value is assigned a color strip. The result is an image consisting of consecutive vertical color strips corresponding to consecutively varying loudness levels.
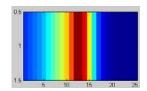


Figure10. Loudness Intensity Rising then Falling

## III. SPEECH VISUALIZATION

### A. Speech Recognition

In case the sound is speech, recognizing speech is of great significance to the hearing impaired. The vocabulary set could be of any size and is dependent on the training. We have used a set composed of 12 words to complete a few sentences with an additional 'Unknown?' word that is displayed whenever a preset threshold is exceeded indicating that either the sound is far from being speech or is not close enough to any of the existing 12 words. Standard preconditioning steps, including pre-emphasis, were performed, and the Cepstrum method was used in identifying word segments. Following the application of the Cepstrum method, the Al-Alaoui Cloning Algorithm [8-9] (equivalently the more recent Boosting method) was applied on the training data (cepstral vectors) for 10 iterations to yield a modified average cepstral vector for each word. Fig. 11 shows the result after each of the 10 iterations. Each of the twelve rows corresponds to a word, and the 6 columns per word correspond to the number of times each word was trained. In each of these matrices, 1's correspond to errors while 0's correspond to correct classifications.



**Figure 11. Error (Misclassification) Reduction using the Cloning Algorithm.**

The resulting word is displayed in a central sub-window



Figure 12. Speech Recognition Display of the word 'Sign'

### B. Speech to ASL Translation

In a deaf community, some people prefer communicating in sign language such as ASL which is common is the United States and Canada. Moreover, people who are initially born deaf or with severe hearing impairments face more difficulties acquiring formal written English language in school as they do not hear it to incidentally acquire it with ease. Hence, it is important to include speech to ASL or Signed English translation within the entire program package. The method of implementation is similar to that of standard speech recognition with the exception that instead of displaying words upon successful matching, static sign pictures are

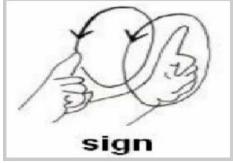displayed consecutively in a concatenated frame-like fashion. Each of the 13 words has its own display image.



Figure 13. Sign Language Display

### C. Single Icon Display

The idea of Single Icon Displays [4] was also implemented. Generally, these displays are limited in that they require prior knowledge of the signal expected, but may be quite useful for certain basic events or if the conditions behind a metric match is kept general. For instance certain visual patterns appearing on the spectrogram may be recognized to cause an icon to pop up indicating 'phone ringing' or 'door knocking', etc. For more relaxed conditions, if the mean pitch value falls in a specific region of the human pitch range and according to the loudness/intensity level, an icon may display such interpreted information accurately. We have successfully implemented 7 icons demonstrated as follows:

TABLE I.    SINGLE ICON DISPLAY

| | |
|---|---|
| 'Possible Male Speaking' |  |
| 'Possible Female Speaking' |  |
| 'Possible Male Shouting' |  |
| 'Possible Female Shouting' |  |
| 'Possibly Non-Human Loud' |  |
| 'Possible Telephone Ringing' |  |
| 'Question' |  |

### D. GUI Controls

An optional FIR Bandpass filter is also included in the Graphic User Interface (GUI) to allow the user to set lower and upper cutoff frequencies as desired. This would of course affect the display of the sub-windows eliminating any unwanted frequency regions in the spectrogram for instance. Other GUI controls involve stopping (freezing), running, and exiting the program at any time. A Help button for almost every sub-window may be pressed and depressed to reveal a brief explanation of the purpose of that respective window. For example, the help button for the Dynamic Circles sub-window has been pressed to reveal the illustrative side-text to the left (Fig. 15).



Figure 14. GUI Controls (BP Filter, Stop, Run, Exit)

### Overall Main Window Display

The Main Window includes all the sub-windows discussed. All these sub-windows are dynamically changing with the continuously incoming sound signal. The Main Window Display is shown in Fig. 15.
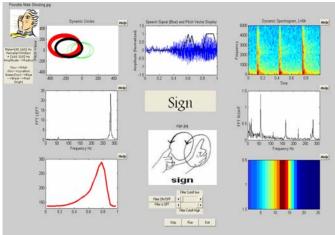


Figure 15. Main Window Display

## IV.    IMPLEMENTATION AND APPLICATIONS

### A. Performance Evaluation

The Program was evaluated by means of a survey conducted in a school for the hearing impaired. The 10 participants had various degrees of hearing impairment from partial deafness to total deafness. The survey focused on assessing the significance and usefulness of the entire program as a unit as well as comparative assessment among the different sub-windows. A summary of the results of the survey is depicted below in Fig. 16 and 17.
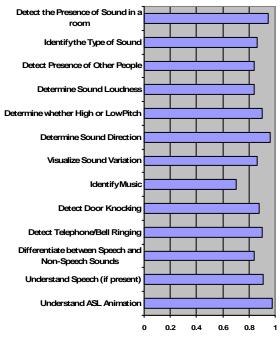
Figure 16. Assessment of functionality

Fig. 16 is based on the participants' assessment of the degree of success of the program in achieving the targets specified.
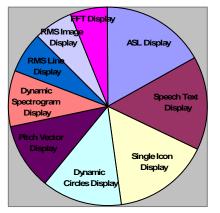


Figure 17. Relative Significance of the Different Displays

The comparative assessment of the various sub-windows (Fig. 17) showed that the participants found that the ASL Display was the most useful followed by Speech Text Display, and that the least useful sub-widow was the FFT Right/Left Display.

The overall usefulness of the entire program prototype was given an average of 88% grade by the participants. The participants were in general excited by the output of the program and showed great interest in the possibility of having this implemented onto a small portable device or cell phone.

### B. Applications

The idea behind the program developed is to provide the hearing impaired with more awareness of the environment around them whether in the workplace or at home. The addition/installation of such software onto any computer or PDA for instance would help reduce communication impediments that are usually faced without disrupting other tasks being performed by the viewer. The entire main window of the program would be displayed for example in a separate corner of the screen while the user would be working in another program. The software could also be integrated into any portable computer based device that would include two built-in microphones or an array of microphones situated on the lower backside of the screen. A small separate device with a reasonable microprocessor supporting the program, sensitive microphones, and a miniature screen may be even used as a portable wrist bound item. The program would continuously provide information of sound source location, loudness, pitch, speech translation, etc.

### V. COMCLUSION AND FUTURE WORK

The program discussed in this paper would provide the hearing impaired with varying types of information extracted from sound. The program was written in the simulation environment of MATLAB. Future work would include actual hardware implementation in more efficient program languages such as C. Speech Recognition may be improved and transformed into a Continuous Speaker-Independent system based on phoneme recognition and HMM. Independent Component Analysis may also be applied on the incoming sound signal. Furthermore, the program may be developed even further to include two-way communication means. For instance, the hardware could be designed to support a two-sided screen; the screen facing the user would show the program display and the screen opposing the user would be used to display any text typed by the user (which simultaneously shows on the screen facing the user). Moreover, with the advances in ASL-to-Speech recognition software, a computer camera may be installed to provide a further option of signing instead of typing text on a keyboard. The sub-windows may be designed with the option of collapsing any frame whenever the user feels it is distracting and/or not required for the moment.

REFERENCES

[1] Ethnologue report for Lebanon
http://www.ethnologue.com/show_country.asp?name=LB

[2] Ross E. Mitchell, " How Many Deaf People Are There in the United States? Estimates From the Survey of Income and Program Participation", Journal of Deaf Studies and Deaf Education, Oxford Journals.
http://jdsde.oxfordjournals.org/cgi/content/full/11/1/112

[3] F. Wai-ling Ho-Ching, Jennifer Mankoff, James A. Landay, (2003). "From Data to Display: the Design and Evaluation of a Peripheral Sound Display for the Deaf", In Proceedings of CHI 2003. 8 pages.

[4] Tarra Matthews, Janette Fong, Jennifer Mankoff, "Visualizing Non-Speech Sounds for the Deaf", In Proceedings of ACM

SIGACCESS conference on Computers and accessibility (ASSETS). Baltimore, MD, pp. 52-59, 2005.

[5] L.R. Rabiner and R.W. Schafer, "Digital Processing of Speech Siqnals", Prentice-Hall, 1978, pp. 150-154.

[6] Sanjit K. Mitra, "Digital Signal Processing A Computer-Based Approach", 3rd ed, McGraw-Hill, pp.533

[7] David Slepian, "Some Comments on Fourier analysis, uncertainty and modeling", SIAM Review, Vol. 25, No.3, pp. 379-393.

[8] Mohamad Adnan Al-Alaoui , "A New Weighted Generalized Inverse Algorithm for Pattern Recognition"; IEEE Transactions on Computers, Vol. C-26, No. 10, pp. 1009-1017, 1977

[9] Mohamad Adnan Al-Alaoui, Rodolphe Mouci, Mohammad M. Mansour, and Rony Ferzli, "A Cloning Approach to Classifier Training", IEEE Transactions on Systems, Man, and Cybernetics—Part A: Systems and Humans, Vol. 32, NO. 6, November 2002, pp. 746-752.

## AUTHORS

**Dr. Mohamad Adnan Al-Alaoui** received his BS degree in mathematics from Eastern Michigan University in 1963, BSEE degree from Wayne State University in 1965, MSEE and PhD degrees in Electrical Engineering from the Georgia Institute of Technology in 1968 and 1974respectively.

His research interests are in Neural Networks and their applications and in Analogue and Digital Signal and Image Processing and their applications in Instrumentation, Communications and Controls. He is currently a Professor at the Electrical & Computer Engineering Department, American University of Beirut, Lebanon. (e-mail: adnan@aub.edu.lb).

**Jimmy Azar** received his BE degree in Electrical Engineering from the American University of Beirut in 2006. His research interests are in Digital Signal and Image Processing, Control Systems, and Medical Imaging. He is currently an MSc student in Biomedical Engineering—Medical Imaging at Technische Universiteit Delft, The Netherlands (e-mail: jimmy.azar@gmail.com).

**Hassan Abou Saleh** received his BE degree in Electrical Engineering from the American University of Beirut in 2006. His research interests are in Wireless Communications, Channel Coding, and Channel Estimation. He is currently an M.A.Sc student in Communications Engineering at Concordia University, Canada. (e-mail: hassan.abousaleh@gmail.com).