

## SOUND SOURCE RECOGNITION FOR INTELLIGENT SURVEILLANCE

Md. Rokunuzzaman<sup>\*1</sup>, Lutfun Nahar Nipa<sup>1</sup>, Tamanna Tasnim Moon<sup>1</sup>, Shafiul Alam<sup>1</sup>

<sup>1</sup>Department of Mechanical Engineering, Rajshahi University of Engineering & Technology (RUET),  
Bangladesh.

<sup>\*</sup>E-mail of the corresponding author: rzaman.mte@ruet.ac.bd

### Abstract

*In surveillance, most of the systems aiming to automatically detect abnormal situations are only based on visual clues while, in some situations, it may be easier to detect a given event using the audio information. A new platform for sustainable development of automatic surveillance is introduced based on sound recognition which gathers information of human behavior, activities and environmental changes. The present research deals with audio events detection in noisy environments for surveillance application. The increasing availability of forensic audio surveillance recordings covering days or weeks of time makes human audition impractical and error prone. The ability of a normal human listener to recognize objects in the environment from only the sounds they produce is extraordinarily robust even in adverse acoustic conditions. In this research, we have developed an intelligent surveillance system which recognizes sound sources and detect events. This system can cover large area which is cost efficient. Sound sources can be recognized by comparing the frequency of sounds. This proposed intelligent surveillance system can recognize different sound sources accurately in real time and pretty much quick.*

**Keywords:** intelligent surveillance; sound recognition; event detection

### 1. Introduction

When the vision system is unable to detect events occurring at a high speed, sound is an important cue for perception. To become intelligent, systems or robots should have to understand situation, make decisions and interact accordingly. Vision system is susceptible to adverse conditions like fog, mist, rain, dark etc. In these conditions sound system can be very effective. If a camera can be made by an intelligent system to respond accordingly for specific sound recognition, an event can be detected instantly. If an intelligent sound system is introduced, it will be easy to detect an event of any unnatural sound for the operator and take actions accordingly. The perspective of using such a recognition system in surveillance and security applications is therefore possible, on the condition that sound class models could be learned and built at the place to control. Long-term audio surveillance recordings may contain speech information and also non-speech sounds such as environmental noise, audible warning and alert signals, footsteps, mechanical sounds, gunshots, and other acoustic information of potential forensic interest. Security system should focus on the robustness of the detection against variable and adverse conditions which is particularly important in surveillance applications. Research in the area of automatic surveillance systems is mainly focused on detecting abnormal events based on the acquired video information [1, 2]. In addition to the traditional video cameras, the use of audio sensors in surveillance and monitoring applications is becoming increasingly important. Audio based surveillance has been studied earlier for detecting various types of acoustic events such as human's coughing in the office environment [3], impulsive sounds like gunshot detection [4], glass breaks, explosions or door alarm [5].

Sound effected camera control is the technology to detect sound sources with the help of the installed sensors in a definite platform and orient the camera accordingly to the direction from where the sound is created or occurred [6, 7]. According to sound events camera movement is controlled and sound source is localized automatically [8, 9]. Our approach to sound classification is inspired by the human auditory system in that we extract auditory features as known from auditory scene analysis from the input signal [10]. At the highest level, all sound recognition systems contain two main modules feature extraction and feature matching [11] Feature extraction plays a very important in the sound recognition process. This is basically a process of dimension reduction or feature reduction as this process eliminates the irrelevant data present in the given input while maintaining important information [12]. The whole process is divided into two stages: training phase and testing phase. Detection is the first step of sound analysis

system and is necessary to extract the significant sounds before initiating the classification step. The classification stage uses a Gaussian Mixture Model classifier with classical acoustical parameters like MFCC [13]. This paper investigates techniques to recognize environmental sounds and their direction, with the purpose of using intelligent techniques in an autonomous mobile surveillance robot.

## 2. System architecture

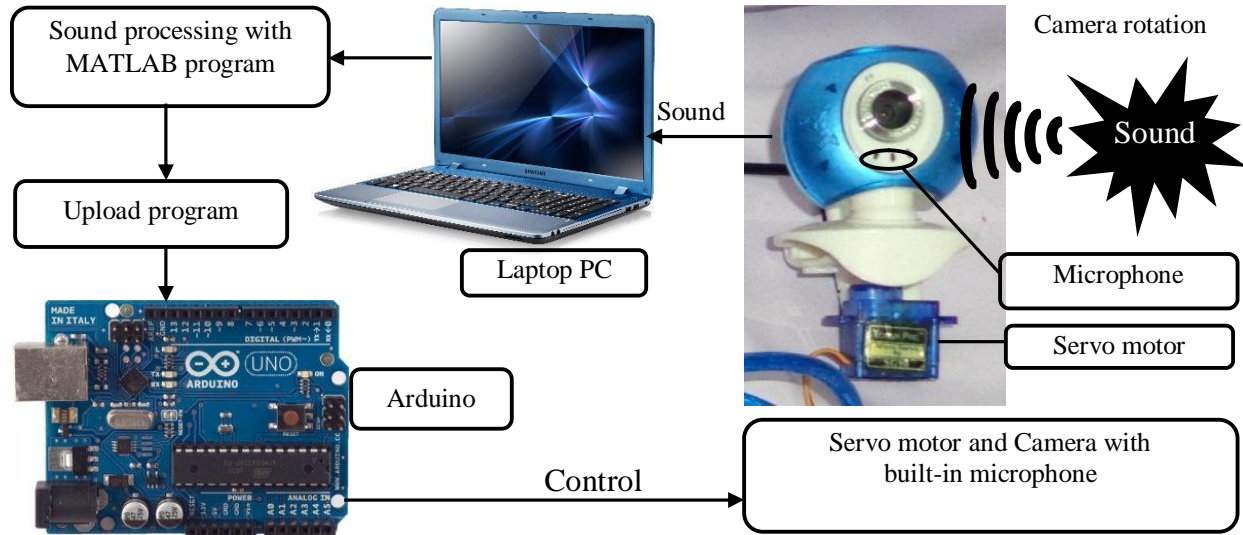


Figure 1: System architecture for development of intelligent surveillance

Figure 1 shows the system architecture for development of intelligent surveillance. When a sound event happens, then it is captured by built-in micro-phone of the camera. The sound signal is passed to the laptop for processing. The sound is processed by the algorithms discussed in section 1 and implemented with MATLAB. The program is then uploaded to an Arduino board. The output signal of the arduino is then fed to the servo motor input to control the movement of the camera toward the sound source.

### 2.1 Configuration & Position of Sound Source

Three sound input are taken for experiments and their position including camera position with specific angle and room layout are shown in Figure 2.

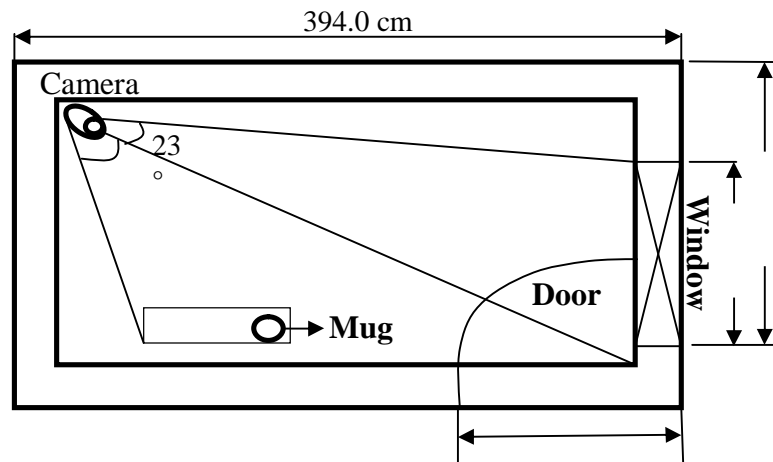


Figure 2: Project room layout

## 2.2 Intelligent surveillance algorithm

The intelligent surveillance algorithm is based on sound recognition and pointing of the camera towards the recognized sound to track objects. The sound recognition is based on feature matching of sound signals between trained signals and input sound. The feature extraction, training and matching is done through a series of operations namely MFCC computation, Vector Quantization and LBG Design algorithms [14]. Figure 3 shows the complete flow chart of the intelligent surveillance algorithm.

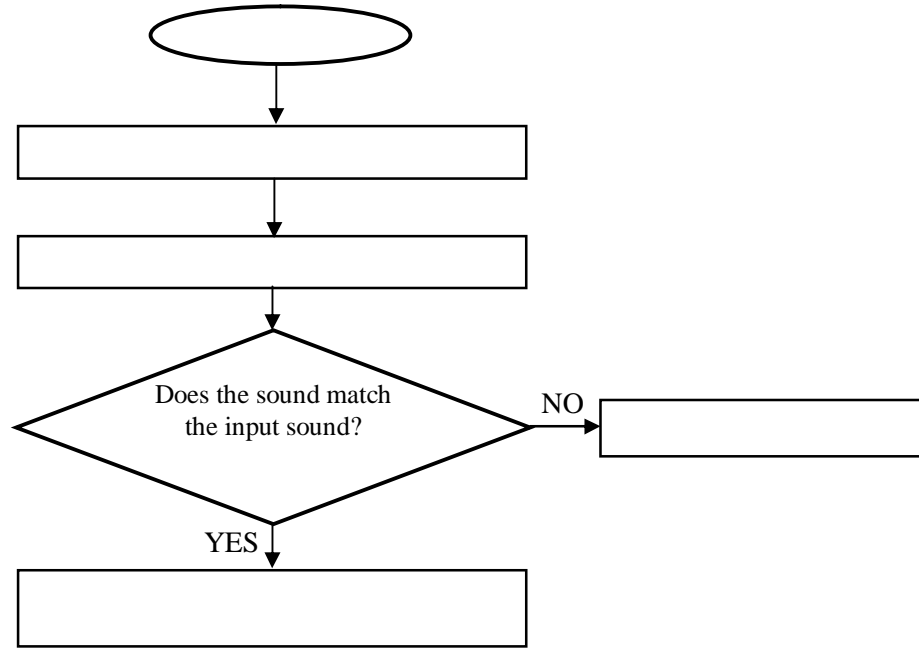


Figure 3: Flow chart of the intelligent surveillance algorithm

## 2.3 Camera movement

After recognition of sound, the recognized signal is transmitted from serial port of PC to Arduino. According to the configuration and position of sound source as shown in figure 2, the Arduino computes the required angle and rotates the camera with the help of a servomotor.

## 3. Results

Different types of sounds from objects have been tested for recognition. The results have been shown in Table1. For each type of sound from object, 5 trials have been taken and success rates have been calculated. The success rate is defined by Equation (1) as

$$\text{Success rate (\%)} = \frac{\text{No. of trials of successful detection}}{\text{Total no. of trials}} \times 100\% \dots \dots \dots (1)$$

Table 1: Success Rates of 3 different sounds

| Sound from Object | Successful detection (Y/N) | Success rate (%) |
|-------------------|----------------------------|------------------|
| Door              | Yes                        | 100              |
|                   | Yes                        |                  |
|                   | Yes                        |                  |
|                   | Yes                        |                  |
|                   | Yes                        |                  |
| Window            | Yes                        | 80               |
|                   | Yes                        |                  |
|                   | Yes                        |                  |
|                   | Yes                        |                  |
|                   | No                         |                  |
| Plastic mug       | Yes                        | 100              |
|                   | Yes                        |                  |
|                   | Yes                        |                  |
|                   | Yes                        |                  |
|                   | Yes                        |                  |

Table 2: Recognition time and Success Rates of different sounds

| Objects     | Average Response Time (s) | Success Rate (100 %) |
|-------------|---------------------------|----------------------|
| Door        | 5                         | 100                  |
| Window      | 5                         | 80                   |
| Plastic Mug | 5                         | 100                  |

The response times of camera in recognizing specific sound source and the success rates of sound source detection of the system are shown in Table 2. Sample sounds of door, window and mug were tested for several times and response time was recorded using stop watch. The obtained average response time is 5 seconds which means it will take 5 seconds to response after the actual sound detection from environment. Success rate of plastic mug is 100% as its sound is different from door and window. As the sound of door and window is similar so success rate of window was not 100%.

Moreover, it is tested that if two sounds occurred at the same time then sound of maximum intensity was detected. Such as when sound of door and mug or window and mug were occurred simultaneously then sound of mug was detected always due to its high intensity. Similarly when sound of door and window were recorded at the same time then door was detected as the sound of door has higher intensity than window.

Table 3: Identification rates of different sounds with code size book

| Code size book | Hamming |
|----------------|---------|
| 1              | 57.14   |
| 2              | 85.7    |
| 4              | 90.47   |
| 8              | 95.24   |

|    |     |
|----|-----|
| 16 | 100 |
| 32 | 100 |
| 64 | 100 |

The identification rates are shown when hamming window is used for framing in a linear frequency scale. Table 3 clearly shows that as codebook size increases, the identification rate for each of the three cases increases when code book size is 16, 32 and 64.

#### 4. Conclusion

Hearing is an important part of normal human interaction, yet we understand surprisingly little about how our brains make sense of sound. This project is driven by the desire to understand how human auditory perception works. This is also to identify the nature of sound and focus the camera directly on the sound source. The time wasted to search the source of the sound will be saved in this project. So this project is more intelligent and efficient from conventional security system. However, the system can be far developed by reducing the response time of camera and more sound can be stored to the system so that it can recognize variety of sounds from the environment. Though the system should have been more robust, the performance of our developed system is quite satisfactory. This project can play an important role in intelligent security system.

#### References:

- [1] Harma, A, McKinney, M. F, Skowronek, J, Automatic surveillance of the acoustic activity in our living environment, *IEEE International Conference on Multimedia and Expo*, July 2005:1-4.
- [2] Michael Cowling. Non- speech environmental sound classification system for autonomous surveillance, *PhD Thesis, Griffith University, Gold Coast Campus*, March 2004
- [3] Clavel, C, Ehrette, T, Richard, G, Events detection for an audio-based surveillance system, *IEEE International Conference on Multimedia and Expo*, July 2005: 1306-1309
- [4]Valenzise, G, Gerosa, L, Tagliasacchi, M, Antonacci, F, Sarti, A, Scream and Gunshot Detection and Localization for Audio-Surveillance Systems, *IEEE Conference on Advanced Video and Signal Based Surveillance*, September 2007: 21-26
- [5] Dufaux, A, Besacier, L, Ansorge, M, Pellandini, F, Automatic sound detection and recognition for noisy environment, *Proc. of the X European Signal Processing Conference*, September, 2000
- [4] Brandstein, M. S, Adcock, J. E, Silverman, H. F, A localization-error-based method for microphone-array design, *Proceedings IEEE International Conference on Acoustics, Speech, and Signal Processing*, May 1996 ; 2: 901- 904
- [6] Wang, H, Chu, P, Voice source localization for automatic camera pointing system in videoconferencing, *IEEE International Conference on Acoustics, Speech, and Signal Processing*, April 1997; 1: 187-190
- [7] Cornaz, C, Hunkeler, U, Velisavljevic, V, An automatic speaker recognition system, *Digital Signal Processing Laboratory*, Federal Institute of Technology, Lausanne, Switzerland, 2003
- [8] Allegro, S, Buchler, M, Launer, S, Automatic sound classification inspired by auditory scene analysis, *Eurospeech*, Aalborg, Denmark, September 2001
- [9] Vacher, M, Istrate, D, Serignat, J. F, Sound detection and classification through transient models using wavelet coefficient trees, *12th European Signal Processing Conference*, September 2004
- [10] Zhong-Xuan Yuan, Bo-Ling Xu, Chong-Zhi Yu, Binary Quantization of Feature Vectors for Robust Text-Independent Speaker Identification, *IEEE Transactions on Speech and Audio Processing*, January 1999; 7(1):70-78

- [11] Gupta, S, Jaafar, J, Fatimah, W, Bansal, A, Feature extraction using MFCC, *Signal & Image Processing: An International Journal (SIPIJ)*, August 2013; 4(4): 101-108
- [12] Soong, F, Rosenberg, E, Juang, B, Rabiner, L, A Vector Quantization Approach to Speaker Recognition, *AT&T Technical Journal*, March/April 1987; 66: 14-26
- [13] Bala A, Kumar A, Birla N, Voice Command Recognition System based on MFCC and DTW, *International Journal of Engineering Science and Technology*, 2010; 2 (12): 7335-734
- [14] Pal, A, Kumar and Sar Anup, An efficient codebook initialization Approach for LBG algorithm, *International Journal of Computer Science, Engineering and Applications (IJCSEA)*, August 2011; 1(4):72-80