# A FRAMEWORK FOR BIOACOUSTICAL SPECIES CLASSIFICATION IN A VERSATILE SERVICE-ORIENTED WIRELESS MESH NETWORK

*Gonzalo Vaca-Castano, Domingo Rodriguez, Julio Castillo, Kejie Lu, Alejandro Rios, Fernando Bird*

AIP Lab/Electrical and Computer Engineering Department, University of Puerto Rico
Stefani Engineering Building, 00680, Mayaguez, Puerto Rico, USA
phone: + (1) 787 832-4040 ext: 2031
email: {gonzalo.vaca,domingo.rodriguez1,kejie.lu}@upr.edu
web: http://walsaip.uprm.edu/

## ABSTRACT

The decline in amphibian populations worldwide has become a tangible example of a major environmental concern due to the fact that amphibian presence has been interpreted as a good indicator of the health of an ecosystem. Consequently, monitoring is an imperative. In this work, a conceptual framework for bioacoustical species classification is formulated and an instantiation of this framework is presented in the form of a sensor array processing (SAP) system. From a digital signal processing perspective, the main feature implemented in the instantiated SAP is an application capability, based on Mel-frequency cepstrum coefficients (MFCC), principal components analysis (PCA), and k-nearest neighbors (k-NN) methods that allows identifying species from audio vocalizations recorded by array sensors. Finally, the processed information is being delivered, through what has been termed a master sensor node (MSN) configuration, to a versatile service-oriented (VESO) wireless mesh network (WMN) which is currently being implemented as an instrumentation testbed at the National Jobos Bay Estuarine Research Reserve (JBNERR) located on the island of Puerto Rico.

## 1. INTRODUCTION

Nearly one-third (32 %) of the world's amphibian species are known to be threatened or extinct according to the Global Amphibian Assessment [1]. Even worse, some estimation suggests that the current extinction rate of amphibians could be 211 times the background amphibian extinction rate [2]. The process of efficiently monitoring amphibian population and correlated changing environmental conditions becomes imperative in light of the decline in amphibian population and the implication of anuran disappearance on overall environmental health. The procedure of finding efficient methods to monitor species and generating inventories is becoming increasingly important as the scientific community struggles to understand reasons behind amphibian declines. Puerto Rico's endemic and critically endangered Puerto Rican crested toad (Peltophryne lemur) is one example of a specie where a declination in population has occurred and recorded. The specie was listed as threatened by the US Fish and Wildlife Service in 1987 (USFWS 1992) and Critically Endangered by the International Union for Conservation of Nature and Natural Resources.

More than 30 national and international organizations are working to save the Puerto Rican Created toad working through the Puerto Rican Crested Toad Species Survival Plan. The efforts include 22 AZA zoos and aquariums which are breeding this species and reintroducing it back to its natural habitat [3]. The decline in amphibian populations in Puerto Rico and worldwide has driven the need to establish more effective monitoring strategies. Wildlife monitoring is challenged by issues such as access to remote sites and limited human resources to deal with strenuous tasks. Automated data recorders (dataloggers) are the most common modern monitoring tools to record amphibian calling activity. Those programmable units are able to record high quality audio during user preset periods. Even though they are more versatile that traditional recording units, they have some inherent problems associated such as the inevitability of weekly field visits to download the information on the loggers, batteries replacement, and null processing capability. The development of new equipment that avoids constants visits to the field is an imperious necessity.

The work described on this paper, present the implementation of a sensor array processing (SAP) unit built using off-the-shelf technology approach. The SAP unit is one the nodes of our purposed Versatile Service Oriented Wireless Sensor Networks(VESO-MESH) which is being currently deployed at Jobos Bay National Estuarine Research Reserve located on the south of the Puerto Rico's island. The infrastructure designed has turned in a bioacoustical framework where different studies are currently being performed. Our main objective with this work is the acoustic Environmental Surveillance Monitoring (ESM) of species. To achieve this objective, we have developed algorithms for automatic identification of species and individuals from their vocalizations which run in the SAP with the purpose of generating alerts about the possible presence of a particular specie in a monitoring process.

## 2. BACKGROUND AND RELATED WORKS

### 2.1 Wireless Sensor Networks

Wireless Sensor Networks (WSN) has become extremely popular in the last years. In fact, in 2002 an already famous survey was published [4], where possible WSN applications were presented, besides of factors influencing sensor network design and different communication architectures. Since then, numerous sensor network applications have been proposed in areas such as indoor/outdoor environmental monitoring, health and wellness monitoring, power monitoring, inventory location monitoring, factory and process automation, seismic and structural monitoring, precision agriculture and military [5]. Many of those applications employ acoustic sensors as a key part in the sensing process.

For example, an experimental counter-sniper system called PinPtr [6] measures the time of arrival of muzzle blasts and shock waves from a shot to detect and locate shooters; an application called ExScal [1] (Extreme Scale Wireless Sensor Networking) provides acoustic intrusion detection using multiple sensor and actuators (magnetometer, a microphone, four passive infrared receivers, a photocell, a sounder, and feedback LEDs.); an in-home application called LISTSENse [7] enables the hearing impaired to be alerted of the audible information in their environment (e.g., smoke alarm and doorbell); a volcanic monitoring application, let to observe 230 eruptions and other volcanic events in northern Ecuador (Volcan Reventador) during 19 days in 2006 [8].

## 2.2 Species Discrimination From Their Vocalizations

Many studies for automatically discriminating species from theirs vocalization have been proposed in recent years. These studies been applied to a wide range of species, including farm animals [9], bats, birds [10], and anurans [11]. Birds have been the preferred target among researchers. Time-Frequency representation is the usual scenario to extract features that allow to obtain a classification of audio samples. FFT spectra, spectrograms, Wigner-Ville distributions (WVDs), Mel-frequecy cepstrum coefficients, and wavelets are examples of common scenarios for time-frequency representations used in automatic identification operations. In 1996, Anderson, Dave, and Margoliash [12] used dynamic time warping (DTW) for automatic recognition of birdsong syllables from continuous recordings. In these studies, syllables were represented by spectrograms and classification was performed by matching the spectrograms to predefined prototypes. It was identified that comparisons operations such as matrix correlations and pattern matching of spectrograms tended to be computationally demanding. They applied these methods to non normalized amplitude vocalizations from two bird species recorded in a low noise environment, achieving high accuracy.

Techniques based on metrics to quantify song similarity have been proposed, as well. McIlraith and Card 1997 [13] conducted research on the recognition of songs of six bird species. In their method, the bird songs were represented with spectral and temporal parameters of the songs. They reduced the complexity of the search space by selecting features exhibiting the greatest discrimination. They then used a neural network for classifying the bird songs. Selouani et al. [14] improved the neural network approach by adding a feedback loop to a multilayer perceptron (MLP) network. Most techniques similar to the one used by McIlraith and Card are derived from human speech recognition techniques.

Under noisy environmental conditions, the actual conditions encountered when performing field recordings, it has been demonstrated that hidden Markov models (HMM) are usually more efficient than knowledge-based recognition methods. Trifa, Kirschel and Taylor [15] were able to distinguish songs from 5 species of antbirds that share a common territory in a rainforest environment in Mexico using HMM techniques. Their experiments show that, with noisy recordings, performance was lower but generally exceeding 90%. Other techniques explored include data-mining[16] , support vector machine [17], machine learning [18], and image processing [19].
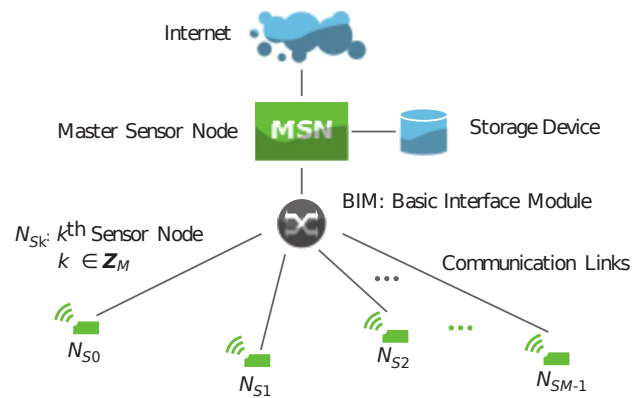


Figure 1: Sensor Array Processing (SAP) Concept

## 3. FRAMEWORK FOR BIOACOUSTICAL SPECIES CLASSIFICATION

The sensor network model depicted in figure 1 correspond to our conceptual representation of our system called Sensor Array Processing (SAP). The SAP is based on a centralized architecture which means there is a central or master sensor node (MSN) that interacts with lower nodes and it is responsible for collecting, processing, and conveying information about the observatory. A sensor array processing (SAP) system, has three essential elements: 1) a set of sensor signal processing nodes (SSP nodes), 2) a wireless routing mechanism (WRM), and 3) a Linux-based high performance embedded computing unit (MSN). The sensor signal processing (SSP) nodes are a set of wireless, low-cost acoustic signal acquisition, storage, and processing nodes which are responsible for sensing the physical signals of interest and they include all the acquisition hardware such as analog to digital converters, audio cards, microphones, etc. These nodes are combined and treated as a sensor array unit without a prescribed topology and their signal-based acoustic information is aggregated through a wireless routing mechanism (WRM) to a Linux-based high-performance embedded computing unit, called a master sensor node (MSN), for further raw data processing and information representation. In this manner, a sensor array processing system or SAPs, each with its own master sensor node unit, may be looked as a node of a complete system forming a Versatile Service-Oriented Wireless Mesh Network for Disaster Relief & Environmental Monitoring, as is illustrated in the figure 2. Note that, several SAP can be included as nodes connected through the MSN, to form part of a "Wireless Bubble". The "wireless bubble" could be connected to the internet if an internet access is available through at least one internet gateway. A versatile service-oriented (VESO) wireless mesh network architecture aims to integrate distributed processes and storage devices into the network so as to provide effective data process and access inside the network, construct high-throughput backbone in Wireless Mesh Networks (WMN) emphasizing the transmission of a larger volume of data, and design and implement a set of service-oriented protocols to efficiently provide services to users inside the network and effectively utilize the network resources in terms of energy, processing, storage, etc. It also expected that can be quickly built to respond to natural disaster, establishes a data intensive environmental monitoring application specifically addressing hurri-
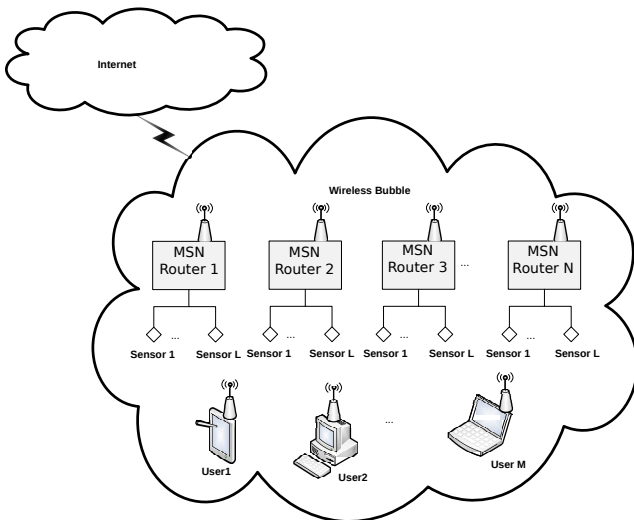
---

Figure 2: VESO Wireless Mesh Network Concept

canes and earthquakes.

## 4. IMPLEMENTATION RESULTS

### 4.1 Equipment

This section provides information about the implementation results pertaining to the SAP framework. A description of the SAP configuration follows. Each SSP node consists of an assembly of three Gumstix[TM] boards. Gumstix[TM] offers ARM based mini board computers with several expansion boards, which can be put together in a single package, attached using the expansion bus included in all the boards. The core of a SSP node is the gumstix verdex XL6P board which contains a processor Marvell PXA270 with XScale[TM] running at 600 Mhz with a memory of 128MB RAM and 32MB Flash. The assembly was completed with two additional boards: the audiostix2 board which provides USB client support and sound card for audio recording and playing; and the netwifimicroSD FCC card which brings Ethernet and 802.11 connectivity and storage through a microSD port. Each SSP is kept inside a Pelican enclosure for protection against water, sun and environment. The verdex XL6P board contains an ARM based microprocessor and enough memory to run a Linux Open Embedded distribution on it. The Linux operative system allows using the device as a small computer with full functionality. The linux shell can be accessed through either a serial port or the network using SSH which is a common network protocol for remote administration of Unix/Linux computers. The linux shell allows using a feature known as crontable or cronjob to schedule programs or services to run, allowing setup a time table to indicate the date and time where the recordings are executed. A small script is created to setup the audio card mixer, input levels, audio format, frequency sampling, number of bits, number of channels and filename of the recording which is adjusted according to the current date and time of the recording.

The protocol 802.11g was chose as communication protocol and the infrastructure mode was selected. In the infrastructure mode, devices communicate with each other communicate through a central place which is known as Access Point (AP). A wireless access point (AP) is required for infrastructure mode wireless networking. To join the WLAN, the AP and all wireless clients must be configured to use the same SSID. The EnGenius EOC-8610 is a waterproof Access Point designed for outdoor conditions which was used in the testbed.

The master sensor node (MSN) unit selected is an embedded PC in mini-ITX form factor from AOpen[TM], model number i945GMt-FSA. The motherboard is equipped as follows:

- CPU: 2.00GHz Intel R Core 2 Duo T7200.
- Memory: Transcend 2GB So-DIMM DDRII 667MHz.
- Hard disk: Hitachi Travelstar 200GB 7200RPM SATA.
- OS: Linux Fedora core 11

Time synchronization between nodes and MSN is achieved using NTP service. NTP server is running in MSN. Periodically, each SSP node queries to MSN for time synchronization. MSN answers with the actual time information which is used by the SSP node to correct its local time. Two shell scripts running from the master sensor node (MSN) coordinate audio transferences from each SSP node to the MSN. The first script copies the audio files from a SSP node to the MSN using SSH commands, and the second script deletes the transferred file from the origin (SSP node) to free storage space in the SSP.

The SAP testbed is solar powered. The system uses a $24V_{DC}$, 216W solar panel and 2 batteries that provide $24V_{DC}$ and storage the electrical charge of the system. An inverter transform the $24V_{DC}$ output to 120VAC. The Master Sensor Node (MSN) is equipped with a real time clock, and a BIOS system that allows turning on the MSN according to the programmed in the BIOS system. Monitoring process is performed during some time period of the day, so the MSN can be kept turn off during some periods to avoid the power consumption and preserving the charge of the batteries.

### 4.2 Algorithms For Automatic Species Identification From Their Vocalizations

MSN processing capacities is not a limitation to develop applications for automatic species identification from audio recordings. In fact, a Matlab application is created with the purpose to generate alerts indicating a possible singing of a particular specie. The objective is to obtain information of available species in real time. Our work pipeline can be divided in 5 different modules. They are: signal preprocessing, segmentation, feature extraction, training and classification.

#### 4.2.1 Signal preprocessing

The quality of the captured audio by recorders in the field is not as good as the records obtained in laboratory controlled conditions. The quality of the records is degraded by several reasons as:

- The quality of the captured audio is lower because the system uses a omni-directional microphone. In an autonomous system, Omni-directional microphones must be used because there is not a user to move the microphone in the direction of the sound origin.
- Environmental factors such as wind and rain introduces low frequency noises.
- The Blur effect caused by multi-path and echoes.
- The distance from the sound source to the recorder has an effect on the amplitude of the signal.

- The object of interest is frequently recorded together with unwanted sound sources such as insects, or animal crowds.

Noise reduction is indispensable because the performance of a classifier dealing with a low SNR ratio audio is poor, but there is not a universal solution to solve the problem. Linear adaptive filters is used to perform noise reduction. The main objective of adaptive filtering is to improve the quality of a signal according to an acceptable criterion of performance using optimum and statistical signal processing techniques that can modify the parameters (coefficients) of a filter during normal operation (usually in real time) without any intervention from the user. The LMS algorithm or other algorithms like Recursive-Least-Square (RLS) can be applied to solve the noise cancellation problem. The purpose of an adaptive noise canceler is to subtract noise from a received signal in an adaptively controlled manner to improve the signal-to-noise ratio. The signal of interest $s(n)$ is corrupted by uncorrelated additive noise $v_1(n)$, and the combined signal $s(n) + v_1(n)$ provides what is known as primary signal. A second reference signal represented by $v_2(n)$, is a signal modeling noise that is uncorrelated with the signal $s(n)$ but correlated with the noise $v_1(n)$. The adaptive noise canceler consists of an adaptive filter that operate in the reference signal output to produce an estimate of the noise $y(n) \approx v_1(n)$,by exploiting the correlation between $v_1(n)$ and $v_2(n)$, to be subtracted from the primary input. Here, the assumption that the signals $s(n)$, $v_1(n)$, and $v_2(n)$ are jointly wide-sense stationary with zero mean values is done. Models of the noise are obtained and classified from audio samples that contain only perturbing noise at different time periods and under different environmental conditions.

### 4.2.2 Segmentation

In the segmentation, raw data is divided into smaller significant objects. Here, those smaller significant objects are the syllables which can be defined as the sound produced by a frog or bird with a single blow of the air from the lungs. Syllables, can be also though as the minimum meaningful unit in a vocalization. the feature extraction is properly applied over each one of the segmented syllables to obtain a particular vector that represents the syllable. The systems that are designed for classifying audio signals usually take segmented audios as initial input because usually it is easier to build analysis and classification systems for segmented objects than for raw data. The segmentation implies a reduction of data and also calculations to be performed in the afterward stages of the identification process. The segmentation stage, presented here, advocates for a segmentation of the vocalization in syllables in an automatically manner.
The segmentation of the syllables is based on obtaining energy function in overlapping frames of audio data and a time-domain analysis similar to the one proposed by Fagerlund [20]. Maximum Energy of the piece of audio being processed, is used to establish the initial threshold. The threshold is a value used to separate a sound caused by specie (i.e. belongs to a syllable) from another caused by background noise. The main idea of the algorithm is consider a group of consecutive frames with a energy value over the threshold value as a syllable. In case, the energy value of the segment is lower than the threshold, it could be considered as noise. the vale of the energy of the noise segments are used to con-

tinually updated the thresholds levels.
Energy levels in the frames can be corrupted by temporal variations of the level signal causing violation of the thresholds levels in an unexpected way. For that reason, merging is applied between groups of segments categorized as syllables which are close, according to factors such as the minimum and maximum syllable duration. Finally, a vocalization is considered finished when the vocalization exceeds a parameter determining the maximum duration of the vocalization or if an inter-syllable gap is prolongating more than the expected inter-syllable gap duration.

### 4.2.3 Feature Extraction

In our framework, features are the set of numbers taken from the data or their transformations that characterize a syllable. Relevant information from the data is obtained by performing a time-frequency domain analysis tool known as the Mel-frequency cepstrum (MFC). The calculation of Mel-frequency cepstral coefficients (MFCCs) starts with the dissection of a signal into frames of 256 pointa which are over-lapped 50%. Each frame is pre-emphasized using a Hamming window, and then a Discrete Fourier Transform (DFT) is executed. Each one of the Fourier coefficients is squared and the result is filtered by a set of 27 Mel-scaled triangular filters. Later, a Discrete Cosine Transform (DCT) is executed transforming the signal to the cepstral domain. From the 27 coefficients, only the first 12 coefficients are hold, because they enclose most of the spectral information. Actually, different authors have revealed that the effect of increasing the number of MFCCs on vowel classification performance is negligible [21]. To finish, a liftering process is accomplished resulting in a smoother signal. The outcome of the preceding operation is a 12 x N feature matrix, where N represents the number of overlapping frames that can be prearranged in the signal corresponding to the syllable. The information enclosed in the feature matrix is vastly redundant. Hence, a Principal Component Analysis (PCA) is performed with the intention of reducing the dimensionality of the feature matrix. After performing the PCA in this set of data, interpreting the feature matrix as a group of 12 vectors of dimension N, we will get N eigenvectors. To reduce dimensionality, only the first 3 of the eigenvectors are chosen. Latterly, a matrix of 12 x 3 elements is achieved which is reorganized as a single vector of 36 dimensions. One bonus benefit of the method is that any syllable is represented by a vector with always the same dimension independently of the length of the syllable.

### 4.2.4 Training

A subset of the total available data is employed as training data. A set of training samples representative of their respective classes is required in order to train the network to perform the pattern recognition on samples not necessarily present in the training set. Each processed syllable produces a feature vector with an associated known class which generates a knowledge base to characterize each one of the classes. A simple interface developed in Matlab was created to add training data in a Postgres$^{TM}$ database which contain besides the information about the feature vector and the classes which each syllable belongs to, information about the file, the sensor, the MSN, and additional features not selected as part of the feature vector which can be useful in data analysis.

### 4.2.5 Classification

The k-nearest neighbor algorithm (k-NN) is selected as classification method. K-nearest neighbor is a supervised learning algorithm where the result of new instance query is classified based on majority of k-nearest neighbor category. The purpose of this algorithm is to classify a new object based on attributes and training samples. The classifiers do not use any model to fit and only based on memory. Given a query point, we find k number of objects or (training points) closest to the query point. The classification is using majority vote among the classification of the k objects. Classification is done based on the distance measure between the test syllable feature vector and the model feature vectors. In our case, the distance or similarity between instances is determined by the euclidean distance and k is set equal to 3.

## 5. CONCLUSIONS

We explored a novel framework for the environmental surveillance monitoring based on a Sensor Array Processing (SAP) module which can be deployed as part of a Versatile Service-Oriented Wireless Mesh Network. Implementation is carried on using off-the-shelf components which are able to provide low cost, moderate power consumption and medium computational capacities. The setup is currently tested on Jobos Bay National Estuarine Research Reserve located on the southern Puerto Rico island.

A major contribution of this work are our algorithms to automatically identify species from the recorded audio in the field by the SAP Unit. SAP capacity allow that the processing can be carried on the device which is located in the field, letting to be integrated to our VESO wireless mesh network.

## REFERENCES

[1] G. A. A. (GAA) tech. rep., The IUCN Red List of Threatened Species, 2008. http://www.iucnredlist.org/amphibians.

[2] M. L. McCallum, "Amphibian decline or extinction? current declines dwarf background extinction rate," *Journal of Herpetology*, vol. 41, no. 3, pp. 483–491, 2007.

[3] "Puerto rican crested toad species survival plan," tech. rep., American Zoo and Aquarium Association (AZA). http://crestedtoadssp.org/.

[4] I. F. Akyildiz, W. Su, Y. Sankarasubramaniam, and E. Cayirci, "Wireless sensor networks: a survey," *Computer Networks*, vol. 38, pp. 393–422, March 2002.

[5] J. Yick, B. Mukherjee, and D. Ghosal, "Wireless sensor network survey," *Computer Networks*, vol. 52, pp. 2292–2330, August 2008.

[6] A. Lédeczi, A. Nádas, P. Völgyesi, G. Balogh, B. Kusy, J. Sallai, G. Pap, S. Dóra, K. Molnár, M. Maróti, and G. Simon, "Countersniper system for urban warfare," *ACM Transactions on Sensor Networks*, vol. 2, pp. 153–177, 2005.

[7] C. R. Baker, K. Armijo, S. Belka, M. Benhabib, V. Bhargava, N. Burkhart, *et al.*, "Wireless sensor networks for home health care," in *Advanced Information Networking and Applications Workshops, 2007, AINAW '07. 21st International Conference on*, vol. 2, pp. 832–837, 2007.

[8] G. Werner-Allen, K. Lorincz, M. Welsh, O. Marcillo, J. Johnson, M. Ruiz, and J. Lees, "Deploying a wireless sensor network on an active volcano," *IEEE Internet Computing*, vol. 10, pp. 18–25, 2006.

[9] G. Manteuffel, B. Puppe, and P. C. Schn, "Vocalization of farm animals as a measure of welfare," *Animal Behaviour Science*, vol. 88, pp. 163–182, 2004.

[10] A. Selin, J. Turunen, and J. T. Tanttu, "Wavelets in recognition of bird sounds," *EURASIP Journal on Advances in Signal Processing*, vol. 2007, p. 9 pages, 2007.

[11] H. Wen, N. Bulusu, C. T. Chou, S. Jha, A. Taylor, and V. N. Tran, "Design and evaluation of a hybrid sensor network for cane toad monitoring," *ACM Trans. Sen. Netw.*, vol. 5, no. 1, pp. 1–28, 2009.

[12] S. E. Anderson, A. S. Dave, and D. Margoliash, "Template-based automatic recognition of birdsong syllables from continuous recordings," *The Journal of the Acoustical Society of America*, vol. 100, pp. 1209–1219, 1996.

[13] A. L. McIlraith and H. C. Card, "Birdsong recognition using backpropagation and multivariate statistics," *IEEE TRANSACTIONS ON SIGNAL PROCESSING*, vol. 45, pp. 2740–2748, November 1997.

[14] S.-A. Selouani, M. Kardouchi, E. Hervet, and D. Roy, "Automatic birdsong recognition based on autoregressive time-delay neural networks," in *Congress onComputational Intelligence Methods and Applications,ICSC 2005*, 2005.

[15] V. Trifa, A. N. G. Kirschel, C. E. Taylor, and E. E. Vallejo, "Automated species recognition of antbirds in a mexican rainforest using hidden markov models," *Journal of the Acoustical Society of America*, vol. 123, pp. 2424–2431, Apr. 2008.

[16] E. Vilches, I. A. Escobar, E. E. Vallejo, and C. E. Taylor, "Data mining applied to acoustic bird species recognition," in *ICPR '06: Proceedings of the 18th International Conference on Pattern Recognition*, (Washington, DC, USA), pp. 400–403, IEEE Computer Society, 2006.

[17] S. Fagerlund, "Bird species recognition using support vector machines," *EURASIP J. Appl. Signal Process.*, no. 1, 2007.

[18] C.-J. Huang, Y.-J. Yang, D.-X. Yang, and Y.-J. Chen, "Frog classification using machine learning techniques," *Expert Systems with Applications*, vol. 36, pp. 3737 – 3743, March 2009.

[19] T. S. Brandes, P. Naskrecki, and H. K. Figueroa, "Using image processing to detect and classify narrow-band cricket and frog calls," *The Journal of the Acoustical Society of America*, vol. 120, pp. 2950–2957, November 2006.

[20] S. Fagerlund, "Automatic recognition of bird species by their sounds," Master's thesis, Helsinski University of Technology, 2004.

[21] H. Mei-Ling Meng, "The use of distintive features for automatic speech recognition," Master's thesis, Massachusetts Institute of Technology, September 1991.