

Developing Birds Sound Recognition System Using an Ontological Approach

Yauheniya Zianouka and Dzianis Bialiauski
Liesia Kajharodava and Aliaksandr Trafimau
Vitalij Chachlou and Juras Hetseвич
*United Institute of Informatics Problems
National Academy of Sciences of Belarus
Minsk, Belarus
{evgeniakacan, dzianis.bialiauski,
lesia.cordell, cncntrt,
vitalikhokhlov, yuras.hetseвич} @gmail.com*

Vadim Zahariev and Kuanysh Zhaksylyk
*Belarusian State University of
Informatics and Radioelectronics
Minsk, Belarus
zahariev@bsuir.by, kuanysh.zhk@gmail.com*

Abstract—The article presents an intelligent model of automated voice recognition systems (on the example of birds). To develop it, a dataset of birds' voices was annotated and processed using Mel-Frequency Cepstral Coefficient as an effective tool for modelling the subjective pitch and frequency content of audio signals. For composing and training the model, Convolutional Neural Network is used to implement high level results. The possibilities of using ontological approaches and OSTIS technology for further improvement of the quality of ML models are shown.

Keywords—recognition system, machine learning, dataset, automatic processing, Mel-frequency cepstral coefficients (MFCCs), Convolutional Neural Network, EfficientNet, ontological approach

I. Introduction

There are many voice signals in nature, and each type of sound signal has its own function. In general, sound signals can be divided into singing and other voices (for example, streams, noise). Singing is a more melodic type of bird voice. It is usually longer and more complex than the stream. Due to many variations of the sound signals of different bird species, there is a problem of their recognition. The relevance of creating such a system is due to the fact that all existing developments on the recognition of animal species are not suitable for Belarusian birds. Only some of them are able to recognize the sound signals of European species, which also affect our project. However, the development and use of an automated system facilitates the recognition process, the system itself has recognition accuracy problems that need to be mentioned. Software for determining the species of animals (for example, birds) by voice signals is based on mathematical calculation models (most often twisted neural networks), which, with insufficient training, can make a computational error that leads to incorrect determination of the biological species we need [1]. Thus, the tasks were set:

- to develop methodological foundations for collecting, annotating and recognizing animal voice signals on the territory of Belarus (in terms of technical implementation).
- to compose a structural scheme for automated recognition of animal voice signals for autonomous continuous monitoring of rare, threatened species and indicator species.
- to increase the accuracy of recognition of animal species (for example, birds).

II. Dataset for training recognition model

A dataset of electronic voice signals corpora is a substantial component for training the recognition model. The primary source of publicly available arrays of animal vocalisations used in the dataset is the Xeno-Canto (<http://www.xeno-canto.org>). Its resources are available for listening, downloading, and studying the characteristics of sound recordings. It is one of the largest sources of audio data of bird vocalisations collected from around the world. The site has API endpoints that can be used to automatically search, download data by scientific or common name of a species or family, region tags, sound types, country, etc.

Machine learning is heavily dependent on data. That's why one of the main tasks that need to be performed during its preparation is the annotation of audio recordings. When developing the method of data annotation for bird voice recognition systems, it is taken into account that the composition of the Belarusian fauna includes rare species of animals and birds, for which there is no or limited scope of annotated sounds and graphic data [2]. It is also considered that the selected data may contain different sounds of the same species, sounds of many other species that are heard together with this species.

In order to improve the results of processing and recognizing birds' voices, each individual audio file is

labelled with the name of its own species. The data downloaded during the collection phase from open sources may include some part of the annotated data, but needs to review and fix the annotation according to the audio event classes. Each audio part of the signal is detected as a silent audio segment or an audio event. Then it is assigned to its appropriate audio subclass based on the nature of its content. Next, all the annotated information with the corresponding timestamps is written to a text file markup. Since the data is mostly recorded in the birds' natural habitats, each recording may also contain some background information, including various noises from other animals or people. In order to adjust and improve the performance of the recognition algorithm, it is recommended to add and annotate recordings in the base for training the algorithm, where there are no bird sounds, but there is a background sound of the environment in which certain species of birds usually exist.

In our work, the system for recognizing the voices of birds was built on fourteen species: Parus Major (Sinica vielikaja), Fringilla Coelebs (Bierascianka), Turdus Philomelos (Drozd-spiavun), Emberiza Citrinella (Strynatka zvyčajnaja), Phylloscopus Collybita, Turdus Merula (Drozd čorny), Sylvia Atricapilla (Lieska-cornahalouka), Luscinia Luscinia (Salaviej uschodni), Acrocephalus Dumetorum (Carotauka sadovaja), Erithacus Rubecula (Malinauka), Loxia Curvirostra (Kryžadziub-jalovik), Phylloscopus Collybita, Turdus Merula (Drozd čorny), Hippolais Icterina (Pierasmieška), Periparus Ater (Sinica-maskouka), Sylvia Communis (Lieska šeraja). Using the API (Application Programming Interface) a dataset was collected (audio recordings) according to the above mentioned bird species. The number of entries was about two hundred for each species. The criteria by which records were selected for training were as follows:

- The duration of audio recordings is more than three seconds and less than 10 minutes.
- Proximity by distance to Belarus (Minsk).

The following information is available for each uploaded audio file: Bird species (Specific epithet); Bird subspecies of the (subspecific epithet); The group to which the species belongs (bird, grasshopper); The name of the species in English; The name of the person who made the audio recording; The country where the sound was recorded; The name of the area where the recording was conducted; Geographical latitude of the place where the recording was made; The type of sound that a bird makes (singing, streaming, etc.); Gender (female or male), etc.

Currently, a total dataset contains 21737 audio recordings with a planned test sample of 4348 audio recordings. On the basis of this corpus, work on improving recognition and optimization algorithms will continue.

III. Data preprocessing

Before starting to create a machine learning model and predicting or classifying, it is necessary to carry out preliminary data processing (Preprocessing) [3]. It is the first and integral step of machine learning, as the quality of the data and the useful information that can be extracted from it, directly affects the learning ability of the model. The main tasks at the data preprocessing stage are:

- Processing of zero values;
- Data normalisation (its transformation to some dimensionless units);
- Control of outliers. These are not errors, but values that abnormally stand out and can distort the model's operation;
- Processing of categorical features;
- The problem of multicollinearity (the presence of a high mutual correlation between two or more independent variables in the model).

The technology stack used to develop the recognition model includes Python programming language. Open Source Python Libraries are Librosa; Tensorflow; NumPy; Keras; Pandas. Development environments are Jupyter and PyCharm. Optimization of algorithms and metrics for building a recognition model allows reducing training time by orders of magnitude. Therefore, it usually makes sense to start building the model using reduced datasets, successively improving the process parameters (see figure 1).

IV. MODELS FOR BUILDING VOICE RECOGNITION SYSTEMS

To date, there are several approaches to the construction of voice recognition systems.

- Recognition models based on spectrograms. The audio signal is converted into a spectrogram - a visual representation of the signal's frequencies over time.
- Models based on the amplitude component of the signal without analysing the frequency characteristics. Very often, recurrent networks RNN, LSTM and GRUs are used as models in this case.
- Models, based on the synthesis of specific useful characteristics of the signal. Among such characteristics, one can distinguish mel-cepstral coefficients (MFCCs), coding with a linear predictor (LPC), gammatone filters (Gammatone filter banks).
- Hybrid models. Usually include several types of models for the synthesis of the recognition model to get the best recognition quality.
- Transfer learning. Training is based on an already pre-trained model on other data where weights are already present. At the same time, the model is further studied on the available data of audio recordings.

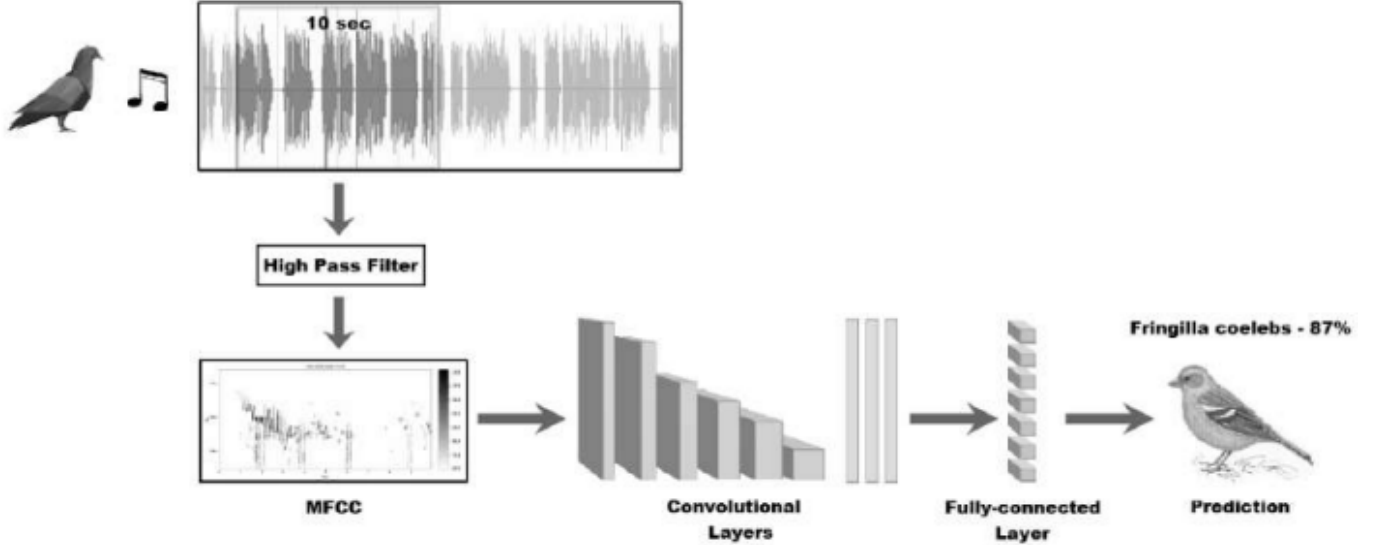


Figure 1: An Audio data preprocessing and neural network model

In our work, we use the first approach, since it certainly results in the best quality of recognized audio signals. For this, we need annotated audio recordings that allow training the model more accurately and recognize bird species.

V. RECOGNITION MODEL BASED ON SPECTROGRAMS

Spectrogram-based models are a powerful tool for sound prediction tasks due to their ability to capture both temporal and spectral features of audio signals. This approach has proven effective in a wide range of applications [4].

Every sound we hear consists of sound frequencies at the same time interval. The essence of a spectrogram is the visualisation of this set of frequencies on a single graph, as opposed to a sonogram, where only the amplitude of the signal is displayed. A spectrogram is a graphical representation of the spectrum of an audio signal as a function of time. It shows which sound frequencies are present in the audio recording at any given time. A spectrogram is built by applying a Fourier transform to short sections of an audio signal called windows. The resulting spectrum is then displayed as a colour map with time on the horizontal axis and frequency on the vertical axis. The colour of each pixel of the spectrogram corresponds to the amplitude of the corresponding frequency.

Mel-spectrogram is a graphical representation of an audio signal in which frequencies are represented on a Mel scale instead of the linear frequency scale used in a conventional spectrogram. The Mel scale is a reproducible scale based on human perception of sound. It is based on the fact that at low frequencies the audio signal can be distinguished with greater resolution, while at higher frequencies

the human ear is less sensitive to changes. Thus, the Mel scale reduces resolution at high frequencies and increases it at low frequencies to better match the human perception of sound. If we combine these two ideas into one, we get a modified spectrogram (MFCC, mel frequency cepstral coefficients), which filters out the frequencies of sounds that a person does not hear, and leaves the most characteristic ones.

Mel scale is calculated as (1):

$$Mel(f) = 2595 \log_{10} \left(1 + \frac{f}{700} \right), \quad (1)$$

$$C_n = \sqrt{\frac{2}{k}} \sum_{k=1}^K (\log S_k) \cos[\pi(k - 0.5)\pi/k], \quad (2)$$

where $Mel(f)$ is the logarithmic scale of the normal frequency scale f . Mel scale has a constant mel-frequency interval, and covers the frequency range of 0 Hz - 20050 Hz. The Mel-Frequency Cepstral Coefficients (MFCCs) are computed from the FFT power coefficients which are filtered by a triangular band pass filter bank. The filter bank consists of 12 triangular filters. The MFCCs are calculated as (2)

where $S_k(k = 1, 2, \dots, K)$ is the output of the filter banks and N is the total number of samples in a 20 ms audio unit.

Since birds sing at high frequencies, a high-pass filter is used to remove unnecessary noise (leave frequencies at a minimum value of 1400 Hz). These coefficients will be sent to the input of the recognition model. The Python library librosa was used to generate the Mel spectrogram: `signal, sr = librosa.load(fp, sr=self.sr, duration=self.duration, mono=self.mono)`

`S_ms = librosa.feature.melspectrogram(y=signal, sr=sr, n_fft=self.n_fft, hop_length=self.hop_length, n_mels=self.n_mels, fmin=self.fmin, htk=self.is_htk,)`