



Urban Sound Quality Recognition Using Smartphones

by

Andreas Stefopoulos

Submitted

in partial fulfilment of the requirements for the degree of

Master of Artificial Intelligence

at the

UNIVERSITY OF PIRAEUS

February, 2023

University of Piraeus, NCSR "Demokritos". All rights reserved.

Author: Andreas Stefopoulos

II-MSc "Artificial Intelligence"

February, 2023

Certified by:

**Mr.Theodoros
Giannakopoulos,
Researcher,
Thesis Supervisor**

Certified by:

**Mr.Maglogiannis Ilias,
Professor, Member of
Examination
Committee**

Certified by:

**Mr.Petasis Georgios,
Professor, Member of
Examination
Committee**

Urban sound quality recognition using smartphones

By

Andreas Stefopoulos

Submitted to the II-MSc “Artificial Intelligence” on February 28, 2023, in partial fulfillment of the requirements for the MSc degree

Abstract

Urban sound quality is a topic that affects the life and well-being of all living creatures, including us humans. Noise pollution in cities is becoming more and more present, resulting in low sound quality levels.

This research project aims to make this phenomenon more visible using artificial intelligence tools and provide insights into how humans' perceptions is functioning under a specific set of sounds. For this purpose, data were collected using a smartphone application and annotated using the perceived sound quality of the user in different locations.

Furthermore, experiments were conducted using audio analysis and neural networks to create a model that can classify the sound quality based on the perceived qualities and provide results and insights based on the user's location.

Lastly, as a final result, an end-to-end process has been made, which could potentially be used every day in different locations, gathering sounds and classifying the quality, aiming to track the change of the quality also in the course of time.

Thesis Supervisor:
Mr. Theodoros Giannakopoulos
Title: Researcher

Acknowledgments

Firstly, I would like to thank my professor Mr. Theodoros Giannakopoulos for providing a concrete approach to my research, guiding me through state-of-the-art technologies to achieve the experimentation, and also for his major contribution to data gathering and support in all topics. Additionally, I am also grateful to the MagCIL lab for supporting also with data gathering.

Contents

1	INTRODUCTION	4
2	BACKGROUND	7
2.1	HUMAN PERCEPTION OF SOUND	7
2.2	AUDIO DEFINITION AND REPRESENTATION	9
2.3	MEL SPECTROGRAM	10
2.4	CONVOLUTIONAL NEURAL NETWORKS	11
2.5	TRANSFER LEARNING	13
3	URBAN SOUND QUALITY AND DATA	16
3.1	URBAN SOUND QUALITY	16
3.2	URBAN SOUND DATA	17
3.2.1	<i>Data Gathering</i>	17
3.2.2	<i>Data Annotation</i>	18
3.2.3	<i>Data Analysis</i>	19
3.2.4	<i>Data Availability</i>	23
4	INFRASTRUCTURE	24
4.1	MOBILE APPLICATION – FRONTEND	25
4.2	MOBILE APPLICATION – BACKEND	30
4.3	DATA VISUALIZATION	31
4.4	DATA PREPARATION - BACKEND	35
5	EXPERIMENTS	36
5.1	DATASETS	37
5.2	CNN MODELS	37
5.2.1	<i>Models Architecture</i>	38
6	RESULTS	40
6.1	RESULTS OF EXPERIMENTATION	40

6.1.1	<i>Model 1</i>	40
6.1.2	<i>Model 2</i>	40
6.1.3	<i>Model 3</i>	41
6.2	TRANSFER LEARNING	41
6.3	FURTHER ANALYSIS	42
6.4	RESULTS INTERPRETATION	43
7	CONCLUSION	44
8	BIBLIOGRAPHY	45

1 Introduction

In recent years, there has been a rise in the amount of focus on the topic of urban sound quality. This is mostly attributable to the influence that a city's acoustic environment may have on the health and quality of life of its residents. Research has shown that exposure to high levels of noise pollution is associated with a variety of health diseases, such as an increased risk of cardiovascular syndrome, disturbed sleep, and stress. Noise pollution has a negative effect on the economy too, in a number of ways, such as the lowering property prices and reducing worker productivity. Therefore, it is crucial to precisely analyze and monitor the sound quality of urban environments to identify places in which interventions may be necessary, in order to lower noise levels and enhance the general livability of cities.

The assessment of sound quality in urban areas has traditionally been implemented using manual methods, such as the deployment of stationary sound level meters at defined places. More recently, however, automated systems have been developed to perform this task. The manual techniques are time consuming and require human labor and could not even offer a complete picture of the acoustic environment at a large scale. These practices depend on the discretion and experience of the people who are collecting the data, where manual procedures can also be prone to mistakes and subjectivity. This is because the judgment and expertise of the people who are doing the gathering, affect the total result.

Nowadays, there has been an increasing interest in utilizing more advanced techniques, such as machine learning and deep learning, to analyze and discover patterns in urban sound data. Specifically, this interest has been fueled by the rise of the big data era. These techniques provide a method of analysis that is both thorough and efficient, and have the potential to considerably increase our capacity to evaluate and monitor the sound quality of urban environments. Algorithms that learn from data gathered and annotated, can evaluate vast volumes of data in a short amount of time and with high levels of precision. They can also recognize patterns and trends that people may have difficulty recognizing. This leads to urban sound quality assessments that are more

accurate, providing a viewpoint on the acoustic environment that is more impartial and neutral.

This dissertation examines the application of machine learning algorithms for the detection of urban sound quality, and the final goal is to come up with recommendations for future research. In order to accomplish this goal, machine learning models are developed, which are capable of classifying urban noises according to the acoustic features of such sounds. A dataset consisting of urban sound recordings that were gathered from a number of different sites is used, in order to evaluate the effectiveness of these models. Through this study machine learning models are developed, that are capable of accurately and reliably recognizing urban sound quality and identifying factors that contribute to variations in sound quality.

In addition to the benefits of urban sound quality assessment and management, the application of machine learning for the recognition of urban sound has the potential to further our understanding of the relationship between the acoustic environment and other aspects of urban life. For instance, the machine learning models that have been trained on urban sound data are able to recognize patterns that are correlated with other factors such as the type of land use, the level of traffic, or the presence of certain types of businesses or facilities. This might also be used as a basis for the creation of interventions that are more focused and successful in the fight against noise pollution and the improvement of urban sound quality.

Furthermore, the use of machine learning for the recognition of urban sounds may potentially have uses outside of the field of urban sound quality. For instance, machine learning models that have been trained on urban sound data may be able to recognize noises that are characteristic of particular kinds of occurrences or activities, such as building construction, car accidents, or public meetings. This has the potential to be useful in a variety of contexts, such as in the areas of public safety, transportation planning, and social media analysis.

The necessity of adequately representing the acoustic properties of urban sounds is one of the primary obstacles that must be overcome when applying machine learning to the task of urban sound detection. In order to do this, it is necessary to implement specific audio processing methods, such as spectrum analysis and feature extraction, in order to

extract the information that is valuable to the investigation. To increase the accuracy and dependability of the models, this process may also entail the creation of bespoke machine-learning algorithms or the use of sophisticated machine-learning techniques, such as deep learning or transfer learning.

The selection and preparation of the dataset utilized to train and assess the machine learning models, is an additional significant part of this project. The dataset accurately reflects the variety and complexity of the urban sound environment, and it contains a sizable number of examples of a wide variety of sounds. The level of quality and resolution of the sound recordings is analyzed during the development, in addition to any potential sources of interference or noise that might have an impact on the precision of the models.

It is crucial to address the practical and operational features of these models and the technical problems involved in developing machine learning models for urban sound quality detection. Additionally, the development of user-friendly interfaces and visualization tools is essential and is implemented, to enable the understanding and analysis of the results, in addition to the integration of the models with other technologies and systems.

In general, the application of machine learning in recognizing urban sound quality provides a rich and difficult research challenge that has the potential to make important contributions to the domains of urban sound quality and machine learning. It is important to evaluate the technological hurdles involved in constructing accurate and efficient machine learning models because the successful development of machine learning models for urban sound detection will considerably increase our capacity to evaluate and regulate the urban sound environment, as well as improve the health and quality of life of people who live in cities.

2 Background

Due to the accessibility of inexpensive resources like computing power and data during the past ten years, artificial intelligence (AI) and machine learning have become nearly universally employed. As a result, AI was able to perform (super)humanly and be developed at a fast pace. In general, most people define AI as “creating machines that are intelligent”. The main drawback of this statement is that it fails to define AI and explain what constitutes an intelligent machine. Although there are many different approaches to the interdisciplinary science of AI, advances in machine learning and deep learning are a paradigm in almost every area of the tech industry.

Although we have heard about sound pollution in our world, only in recent times has AI been used to provide deeper analysis and insights to suppress this phenomenon and identify the root cause of its existence. The definition of sound is when air molecules collide with one another as a result of an object's vibration, sound is produced. These air molecules oscillate, which induces small pressure differences within them that function as sound waves. These waves, which are additionally referred to as mechanical waves, move through a medium while transmitting energy from one place to another. In a deeper analysis, this is the exact reason why sound cannot travel into the vacuum of space, there is simply no medium for it to do so. The period of a wave is the length of time it takes to complete one cycle. The opposite of period is frequency, which is measured in cycles per second and expressed in Hz. The frequency increases and decreases depending on how long it takes a cycle to complete. A wave with peaks that are visually closer to one another would have a greater frequency than a wave with peaks that are farther apart.

2.1 Human Perception of Sound

Regarding human interaction with sound, the pitch of a sound is frequently used to reflect how we perceive frequency. While we use the term "pitch" to describe sounds, frequency refers to the numerical representation of the rate of cyclical repetition of a waveform. The pitch of a sound increases with frequency, conversely, a sound's pitch decreases with frequency. The human auditory system is sensitive to sound frequencies

ranging from 20 Hz to 20,000 Hz, or roughly 10 octaves, which humans hear along the pitch dimension. The basilar membrane traverses the length of the cochlea and vibrates in response to noises that reach the cochlea via the eardrum and middle ear bones' motions. The basilar membrane's operation can be likened to that of a prism in that a typical sound's range of frequencies is scattered to different areas inside the cochlea along the basilar membrane (Figure 1). Every place along the basilar membrane reacts best to a certain frequency, which is referred to as the best frequency or characteristic frequency (CF). In this approach, the frequency content of a sound is encoded in a frequency-to-place map over the length of the basilar membrane, enabling tonotopic organization with a gradient from the variety of frequencies humans can perceive, from the apex to the base of the cochlea. This arrangement is maintained from the cochlea to the primary auditory cortex through the inner hair cells and the auditory nerve, as well as the brainstem and midbrain. It is a fundamental organizing concept for both brain coding and perception.

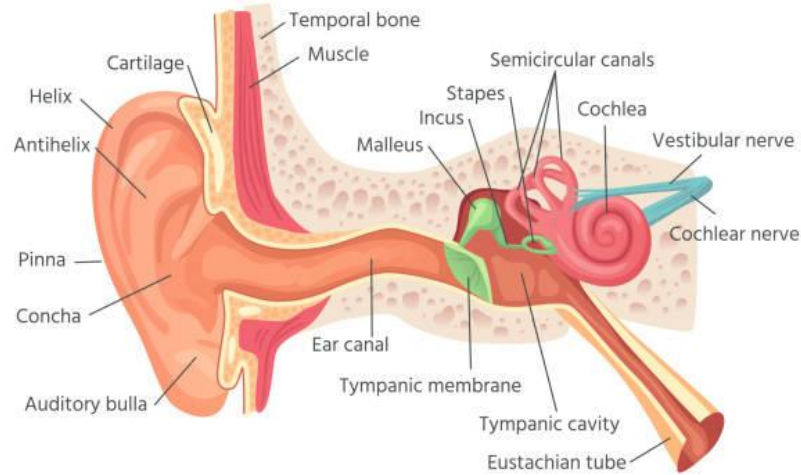


Figure 1: Human Acoustic System

2.2 Audio Definition and Representation

All audio is, by definition, an analog signal. Analog signals are a continuous graph of sound amplitude versus time with infinite values at each infinitesimal unit of time. Storing a raw analog signal would be nearly impossible because it would necessitate infinite storage. Instead, we use a series of operations to extract values from an analog signal at regular intervals. This enables us to store the signals in a digital format while collecting enough data to reproduce a sound in a fraction of the memory. Analog to Digital Conversion (ADC) is a process that uses sampling and quantization to collect a finite set of values for any given analog signal.

Sampling is the process of extracting values at fixed, equidistant time intervals rather than collecting every value in the continuous analog signal. The most common audio sampling rate is 44.1 kHz, which corresponds to 44,100 values per second of sound. This sampling rate allows us to extract all data values within the human hearing range in the most optimal way.

Unlike sampling, which extracts values at fixed time intervals along the horizontal axis, quantization divides values on a waveform's vertical axis into a range of fixed equidistant values. Quantization rounds the exact value at a given time interval to the nearest quantized value when selecting a value at a given time interval. The resolution, or the number of quantized values, is measured in bits. A standard CD has a bit depth, or resolution, of 16 bits, which translates to 65,500 quantized values. The greater the dynamic range when converting an analog signal to a digital signal, the greater the bit depth during quantization.

When a microphone picks up sound, its internal diaphragm oscillates, creating an analog signal that is transmitted to a sound card. The newly created digital signal is sent to the computer for processing by this sound card's ADC.

To manipulate the audio data in Python, the file is transformed into a NumPy array when it is loaded. The amplitude at each timestep of audio is represented in memory as a time series of integers. For instance, a one-second audio clip would include 16800 numbers if the sampling rate was 16800. The data only includes the amplitude numbers and not the time values because the measurements are made at specific intervals of time.

We can determine the time instant at which each measurement of an amplitude number was taken by using the sample rate.

Furthermore, to analyze the sound and process it using AI, we need to generate Mel Spectrograms. This approach is the most optimal in audio processing for classification. The theoretical part behind the creation of a Mel Spectrogram is described in detail in the following chapter.

2.3 Mel Spectrogram

A single-frequency sound wave makes up an audio signal. When we take samples of the signal over time, we only record the resulting amplitudes. The Fourier transform is a mathematical formula that allows us to decompose a signal into its frequencies as well as the amplitude of each frequency. In other words, it converts the signal from time to frequency. The result is referred to as a spectrum. This is possible because each signal can be decomposed into a series of sine and cosine waves that sum to the original signal. This is known as Fourier's theorem, which is a remarkable theorem and is widely used in many sciences. The fast Fourier transform (FFT) is an algorithm that computes the Fourier transform quickly and is commonly used in signal processing. The fast Fourier transform is a powerful tool for analyzing the frequency content of a signal, but in the case of music and speech, the frequency content of the signal varies over time. These signals are referred to as non-periodic signals. Although, we require a method to represent the spectrum of these signals as it changes over time. Thus, we compute several spectrums by performing FFT on several windows of the signal known as the short-time Fourier transform.

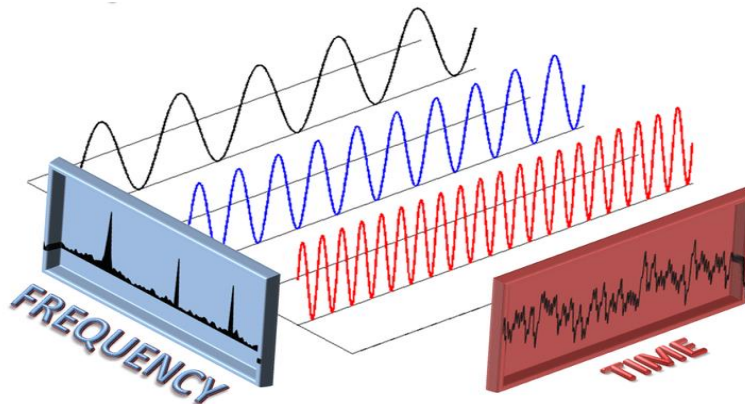


Figure 2: Fast Fourier Transformation

The spectrogram is obtained by performing the FFT on overlapping windowed segments of the signal. A spectrogram can be thought of as a collection of FFTs stacked on top of each other. It is a method of visually representing the loudness, or amplitude, of a signal as it varies over time at different frequencies. When computing the spectrogram, some additional details are taking place behind the scenes. The color dimension is converted to decibels, and the y-axis is converted to a log scale (representing the log scale of the amplitude). This is due to the fact that humans can only perceive a narrow and limited range of frequencies and amplitudes. In Figure 2 we can see the visual representation of FFT.

In the current research, Mel Spectrograms were used in the sound analysis, because humans do not perceive frequencies on a linear scale, according to studies. We detect variations in lower-frequency regions better than in higher frequencies. For instance, we could easily distinguish the difference between 500 and 1000 Hz, but not between 10,000 and 10,500 Hz, despite the fact that the distance between both pairs is the same. Stevens, Volkman, and Newmann recommended a pitch unit in 1937 that sounded equally distant to the listener. The mel scale is used to describe this. To convert frequencies to the mel scale, we use arithmetic and logical operations. In conclusion, a mel spectrogram is a spectrogram where the frequencies are converted to the mel scale and uses the Decibel Scale instead of Amplitude to indicate colors.

2.4 Convolutional Neural Networks

Convolutional neural networks (CNNs) are a type of neural network specifically designed for working with data that has a grid-like structure, such as an image. CNNs are particularly effective for image recognition tasks and have been successful in a wide range of applications, including object classification, object detection, and image segmentation.

The basic structure of a CNN consists of an input layer, an output layer, and one or more hidden layers in between. The input layer of a CNN is typically a multi-channel image, where each channel represents a different feature of the input data (e.g. red, green, and blue channels for an RGB image). The output layer is a vector of class probabilities, representing the predicted class label for the input data.

The hidden layers of a CNN are made of multiple units called "neurons," which are connected to each other through a series of "weights" and "biases." Each neuron receives input from a number of other neurons from the previous layer and then applies a mathematical operation (such as a dot product or a nonlinear activation function) to the inputs in order to produce an output. Activation functions are used to introduce nonlinearity into the model, which is important for learning more complex relationships in the data. Commonly used activation functions include the sigmoid function, the tanh function, and the ReLU (Rectified Linear Unit) function. The output of each neuron is then passed on to the next layer, where it becomes the input for the neurons in that layer.

One key feature of CNNs is the use of "convolutional layers," which are designed to automatically learn a set of filters that can be applied to the input data in order to extract useful features. These filters are typically small (e.g., 3x3 or 5x5 pixels), and are applied to the input data by sliding them across the input grid, applying the dot product between the weights of the filter and the values of the input pixels at each location. This process is called "convolution," and it allows the CNN to learn features that are invariant to translations, rotations, and other geometric transformations of the input data.

Another important feature of CNNs is the use of "pooling layers," which are used to reduce the size of the input data by down-sampling it. This is typically done by applying a max or average pooling operation over a small region of the input data, which effectively reduces the resolution of the input while retaining the most important features. Pooling layers are useful because they allow the CNN to be more robust to small variations in the input data (e.g. due to translation or rotation), and also help to reduce the number of parameters in the model, which can improve generalization performance.

To use a CNN to analyze a spectrogram, the first step is to convert the audio signal into a spectrogram representation. This operation involves applying the Fourier transform to a windowed segment of the signal and then displaying the resulting complex spectrum as a function of time. The resulting spectrogram can then be passed to the CNN as an input tensor, with the x-axis representing time and the y-axis representing frequency.

The CNN can then apply convolutional filters to the spectrogram in order to extract features that are relevant for a given task (e.g., identifying the presence of certain pitches or timbres). CNNs require substantially less pre-processing than other classification

techniques. While approaches used in the past generate filters that are manual-engineered, CNNs can learn these filters/characteristics via training. Through the use of appropriate filters by applying a mathematical approach (multiplication, division) to the original array table with the convolution table, a CNN may successfully capture the Spatial and Temporal correlations in a picture. Due to the dimensionality reduction and the reusability of weights, the architecture performs superior fitting to the picture dataset.

The CNN's function is to compress the pictures into a format that is easier to process while retaining elements that are important for generating a prediction. This is critical for designing an architecture that is not just effective at learning features but also scalable to large datasets. CNNs do not have to be confined to a single Convolutional Layer. In most cases, the first ConvLayer is in charge of collecting Low-Level information such as edges, color, gradient direction, and so on. With further layers, the architecture adjusts to the High-Level characteristics as well, giving us a network that understands the photos in the dataset almost as we would.

2.5 Transfer Learning

Transfer learning is a machine learning technique in which a model trained on one task is re-used or fine-tuned for a second, related task. The idea is to leverage the knowledge and features learned from the first task to improve the performance of the second task. In transfer learning, a pre-trained model is used as a starting point and then further trained on the new data for the target task. This can be useful when there is a lack of labeled data or computational resources for training a model from scratch on the target task. Transfer learning can be applied in a variety of situations, such as natural language processing, computer vision, and speech recognition. It can also be applied at different levels, such as at the level of individual neurons or layers, or at the level of the entire model.

One common type of model used in transfer learning for image classification tasks is a convolutional neural network (CNN). The convolutional layers of a CNN typically consist of a set of filters, which are learned during training and applied to the input image to extract features. The filters are applied at different locations and scales as the input image

is processed through the convolutional layers, allowing the CNN to learn a hierarchy of features at different levels of abstraction.

The fully connected layers of a CNN take the output of the convolutional layers and use it to make a prediction about the class of the input image. These layers are called "fully connected" because each neuron in a fully connected layer is connected to every neuron in the previous layer. The fully connected layers use the features learned by the convolutional layers to make a final prediction about the class of the input image.

To fine-tune a pre-trained CNN for a target task, it is common to unfreeze a few of the final layers of the network and retrain them using the new data for the target task. This allows the model to learn task-specific features that are relevant to the target task, while still leveraging the lower-level features learned from the pre-trained model. This can be done by using a smaller learning rate for the pre-trained layers, which helps to preserve the knowledge learned during the initial training phase. It is also common to add additional layers to the CNN for the target task, which can help the model learn more task-specific features. During the fine-tuning process, it is important to choose an appropriate loss function and optimization algorithm for the target task, as well as to carefully select the hyperparameters of the model.

In other words, transfer learning leverages past assignment expertise to improve prediction about a new task. During transfer learning, the information of a previously trained machine learning model is transferred to a separate but closely related task. If there is a trained basic classifier to predict if a picture contains a specific object, you might utilize the model's training experience to identify additional items. Neural networks in computer vision often seek to identify edges in the first layer, shapes in the middle layer, and task-specific properties in the latter layers. Transfer learning is used in the early and central layers, but the subsequent layers are just retrained. Transfer learning may build an effective machine learning model with relatively minimal training data since the model has previously been pre-trained.

Lastly, training time is reduced because constructing a deep neural network from the beginning of a challenging task might take days or even weeks. When there is not a sufficient amount of annotated data to train our model and a pre-trained model has been trained on similar data and tasks, the best practice is to use the transfer learning model.

Someone could simply restore it and retrain certain layers to focus on the specific case. Transfer learning, on the other hand, works only if the features learned in the first task are universal, in the sense that they may be applied to another activity. In addition, the model's input must be the same size as when it was originally trained, which in the specific case the size of spectrograms is the same for both the pre-trained and final model.

3 Urban Sound Quality and Data

Cities are vibrant and dynamic environments that are frequently accompanied by levels of unwanted noise, which may add to feelings of discomfort or displeasure in many metropolitan regions. The lack of noise and the presence of pleasant sounds can both contribute to feelings of tranquility, relaxation, and “quietness”. This sensation may be encountered in a variety of urban settings, for instance at home, at a park, or on the plaza. However, due to fast and ongoing urbanization, the spatial variance between desired and undesired noises is reducing and may even disappear in some regions. As a result, the features of metropolitan regions where individuals might temporarily retreat from urban stresses like noise may alter or become scarcer.

3.1 Urban Sound Quality

Urban Sound Quality is based on human perception. Apart from perception, there are some general rules applied like the correlation between the volume of sound with the quality. That being said, humans often classify a sound as a “bad” quality based on how loud this sound is. In the current research, there will be a detailed analysis of how different the perception of individuals is and how this is being portrayed in the data gathering process.

In urban areas, there is a variety of sounds, consisting of nature-generated sounds to human-generated sounds. The nature-generated sounds may vary from physical phenomena, like heavy rain, thunderstorms, and airwaves which cause turbulences of human-made objects and constructions, to animal-generated sounds. The human-made sounds vary a lot more, from having sounds of speech, which in crowded areas can be uncomfortable, to machinery and automotive sounds. For this purpose, in order to classify a sound, there has been created a scale from 1 to 5 (and from 1 to 3). In class 1 we label the quality of sound that is unbearable for the user collecting the data and in class 5 as the optimal sound quality. As mentioned, there is a deviation from what we perceive as an optimal sound and also how this scale is applied to sounds from different individuals. A

lot of people living in urban areas may have become familiar with sound pollution, therefore the classification differs in comparison with someone living in a quieter area.

3.2 Urban Sound Data

3.2.1 Data Gathering

For the research of Urban Sound Quality, data was gathered using a mobile application, developed for this purpose. Users in Athens, Greece and Prague, Czech Republic used this application to record a ten-second sound every time and annotate it based on the perceived sound quality. The sounds were gathered in urban areas and at various times of the day. The collection was performed by walking on the streets, in areas with a lack of vehicles, in crowded areas, and in less populated areas in order to gather diverse sounds of either human or nature generated. In Figure 3 we can observe the distribution and the annotated value of the sound in Athens, Greece whereas in Figure 4 we can see the same information for Prague, Czech Republic.

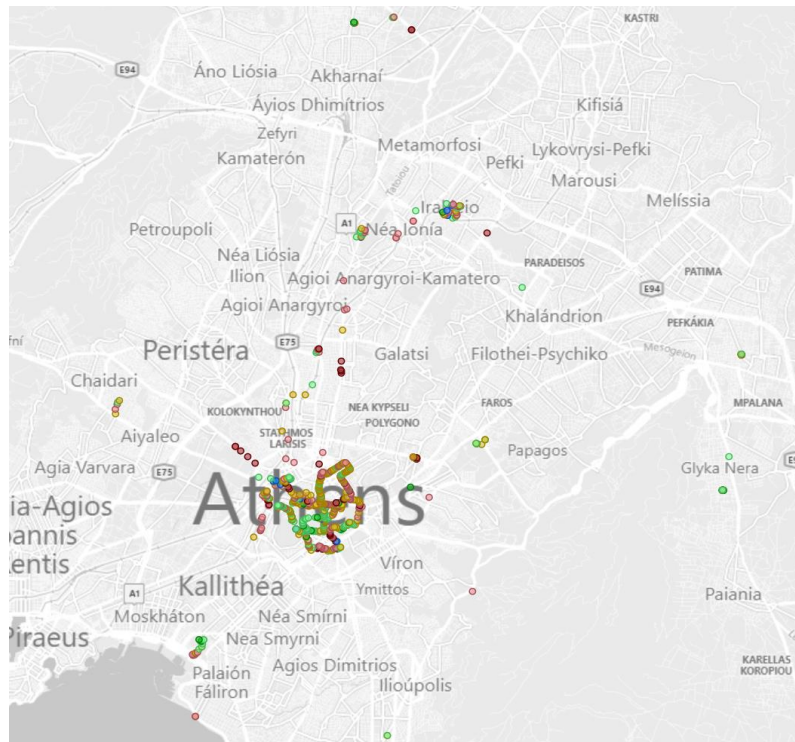


Figure 3: Distribution of sounds in Athens, Greece

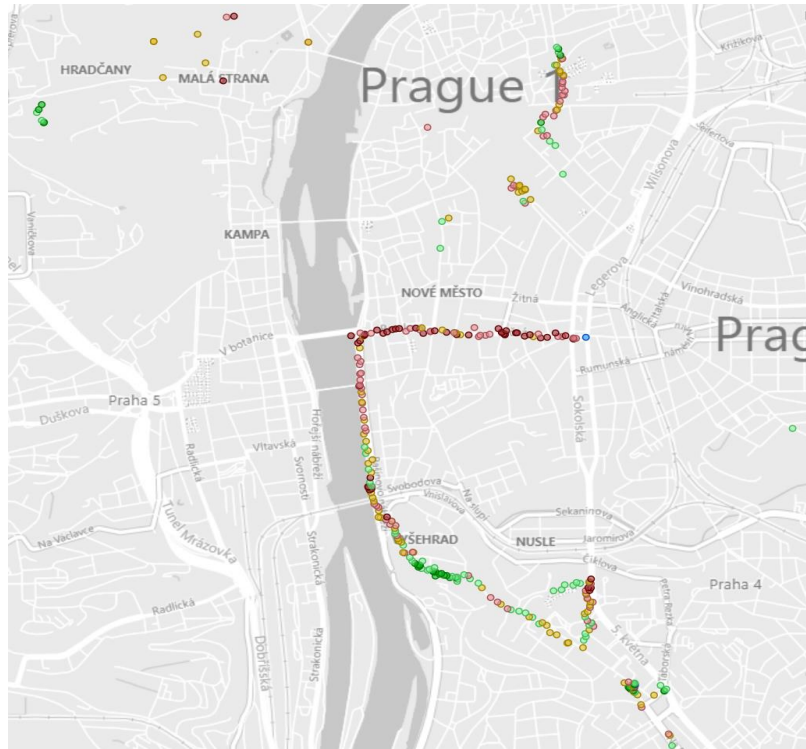


Figure 4: Distribution of sounds in Prague, Czech Republic

3.2.2 Data Annotation

While gathering the data, the user had to annotate the data based on the perceived sound quality on a scale from 1 being unbearable sound, to 5 being optimal sound quality. This was based purely on perception and the only information provided to the users were the aforementioned instructions. Furthermore, the annotation was done right after the recording of the sound, where the user has also the ability to listen to the recording in case of an issue in the process. A common question from the users was that the record is 10 seconds and if in a minor part of the record, there was an unbearable sound and the rest was silence, how would someone classify it? The instruction given in that case is that anything included in the sound is to be reflected in the annotation, meaning that in the case of something being recorded that pollutes the optimal sound, must be appropriately annotated. In some cases, due to the high sensitivity of the new smartphones' microphones, steps can be heard in the recordings. This could not be perceived during the process of recording but could be heard during listening to the recording before sending. In this case, the user's annotation was not affected due to purely depicting the perceived soundscape quality. Additional sounds can also be heard because in most of the cases the

recordings were done while the users were walking, resulting in clothes and breathing sounds apart from steps. The same disciplines apply in these scenarios.

3.2.3 Data Analysis

While gathering data, the piece of information gathered from the user is the Device Identification Number, the Timestamp of the recording in DateTime format, the raw file of the sound record, the longitude and latitude of the user, and the annotated value from the dropdown menu. The data gathered in total per Device ID can be seen in Figure 5.

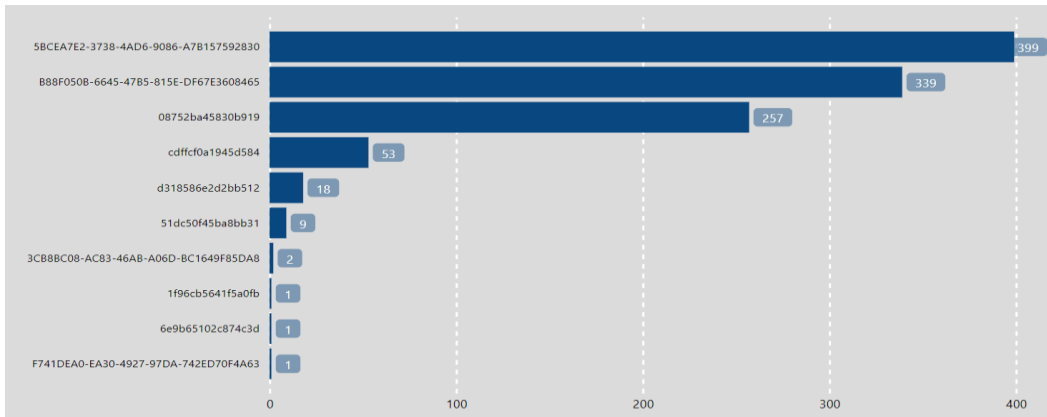


Figure 5: Total Recordings per Device ID

The total number of recordings is 1.080 also having additional 29 recordings from the first version of the application which did not include the Device ID. The average perceived quality per Device ID can be seen in Figure 6.

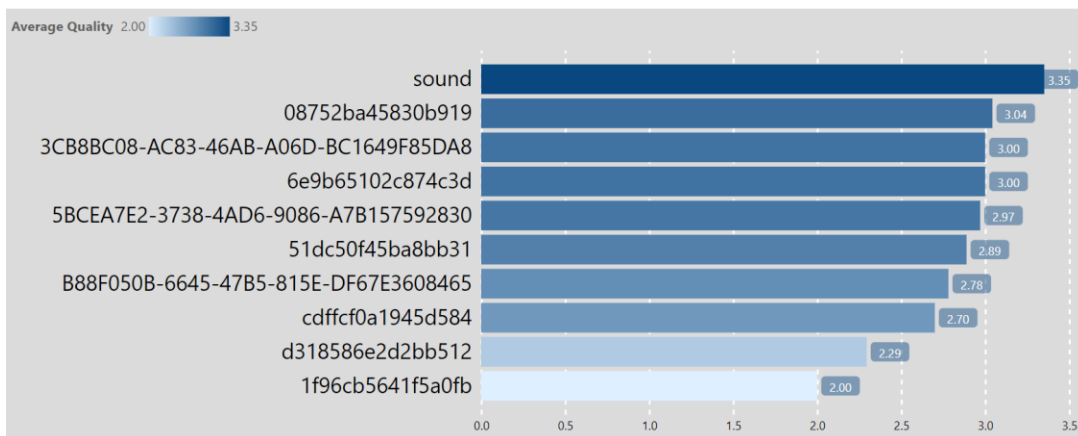


Figure 6: Average Perceived Sound Quality per Device ID

Additionally, the distribution of the perceived soundscape quality in the data can be seen in Figure 7, including the records without Device ID.

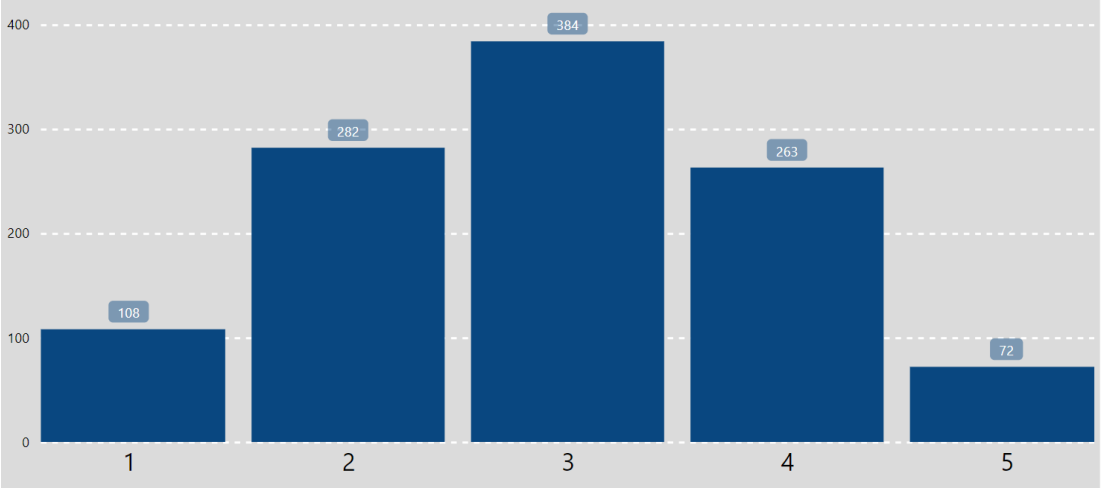


Figure 7: Distribution of Perceived Qualities

The time of day the data was collected also varies, and the distribution can be seen in Figure 8. This variation leads to a wider spectrum of data.

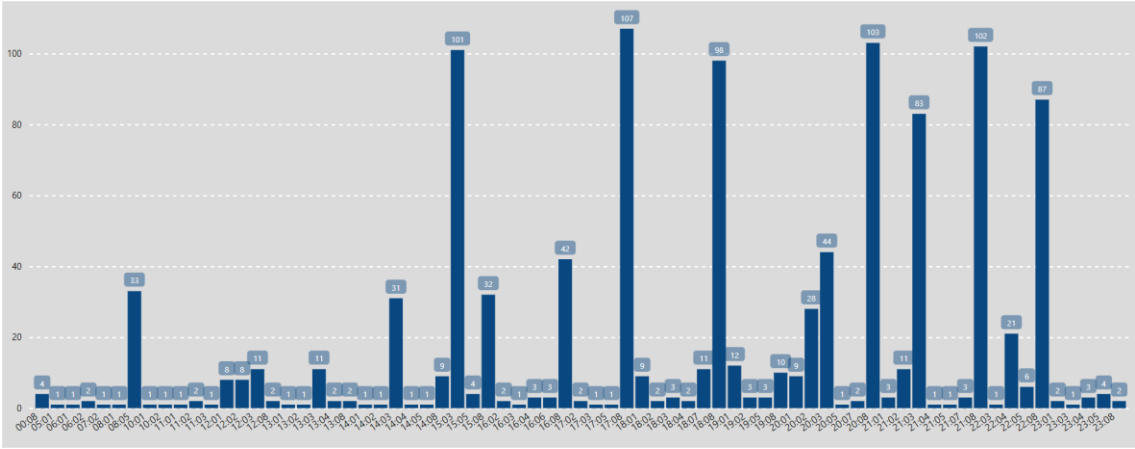


Figure 8: Distribution of Number of Sounds in Time of Day

Average quality can also be investigated in relation to the dimension of time. In Figure 9 we can observe when were the optimal timeslots of higher perceived quality in all of the

samples. The time duration analysis would also provide the correlation between time and quality, as long as the sounds were gathered in the same place, which in our case were not.

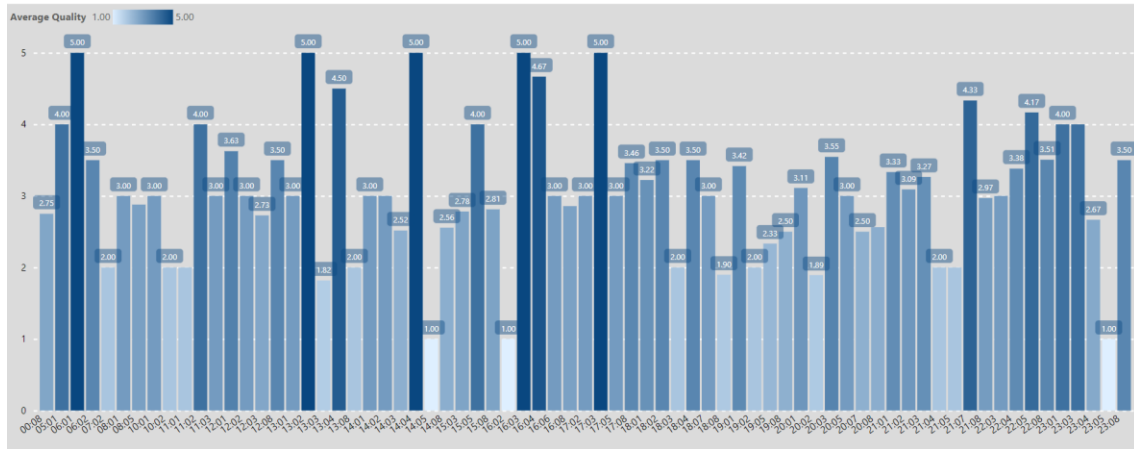


Figure 9: Average Perceived Sound Quality in Time of Day

Gathering data in two different countries is also another dimension worth investigating, where further analysis can provide insights into the differences between the two countries. Starting with information in Figure 9, split by country in Figure 10 and Figure 11.

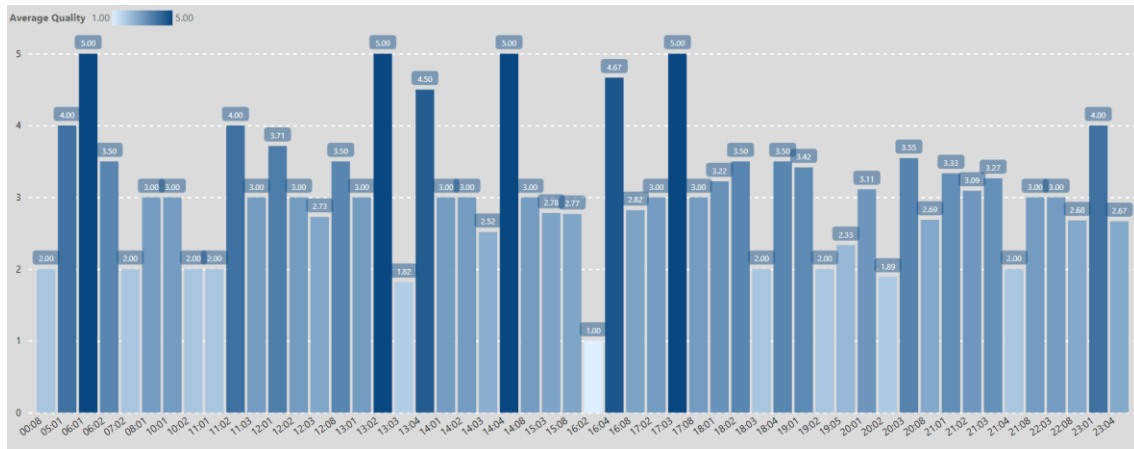


Figure 10: Average Perceived Sound Quality in Time of Day in Athens, Greece

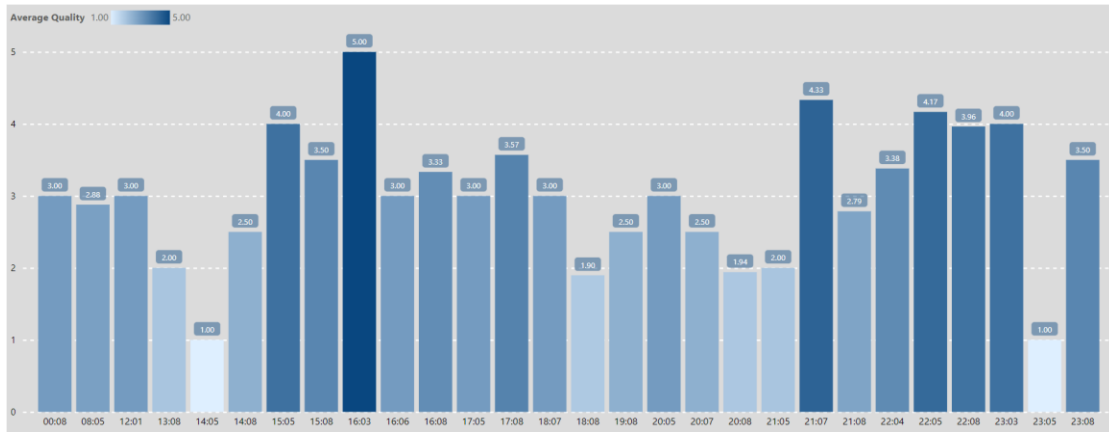


Figure 11: Average Perceived Sound Quality in Time of Day in Prague, Czech Republic

Lastly, in Figures 12 and 13 we have the analysis of the information in Figure 7, the distribution of the perceived soundscape quality per country.

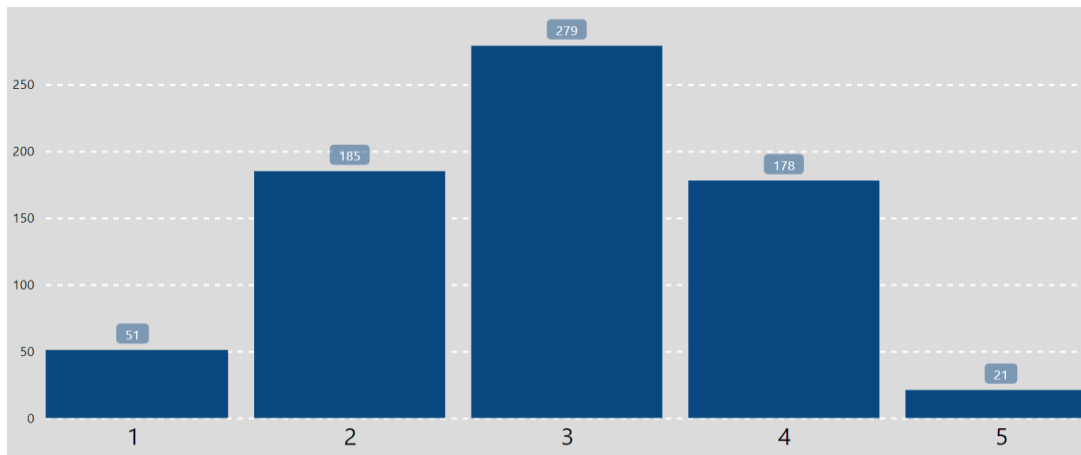


Figure 12: Distribution of Perceived Qualities in Athens, Greece

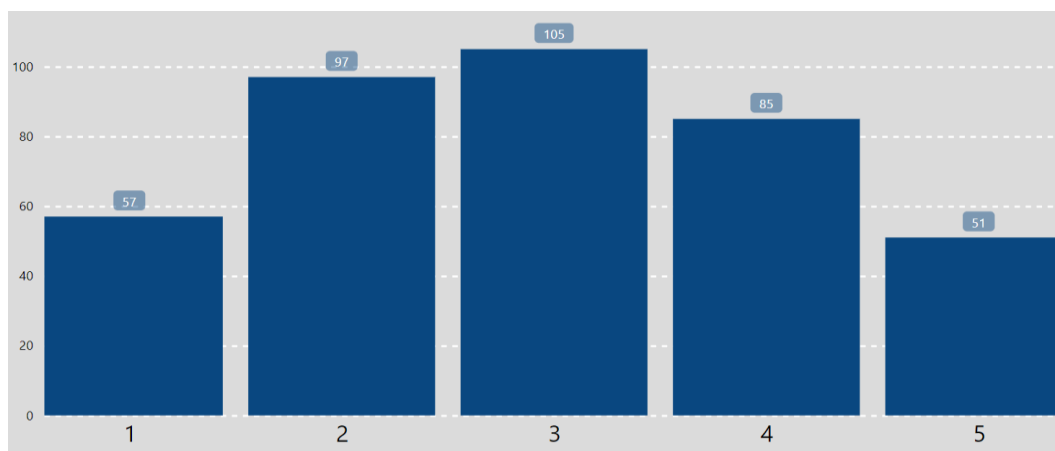


Figure 13: Distribution of Perceived Qualities in Prague, Czech Republic

3.2.4 Data Availability

The recordings gathered using the mobile application can be found in the following repository: <https://github.com/magcil/soundscape-quality/>

The data in the repository are split into train (80%) and test (20%) and the distribution per class is the following:

- train:
 - class 1: 86
 - class 2: 220
 - class 3: 294
 - class 4: 200
 - class 5: 53
- test:
 - class 1: 22
 - class 2: 55
 - class 3: 74
 - class 4: 51
 - class 5: 14

In the data directory, there is also a file “database.csv” where additional information can be found for each recording. The format of the name of each file is: `channel_class_deviceid_longitude_latitude_timestamp`.

4 Infrastructure

For the purpose of the current research, a mobile application was developed which functioned as a data-gathering tool, assisting in the accuracy of the data by relying only on custom-developed code allowing to modify every detail for faster and more concrete results. The mobile application was developed using React Native and Node JS. The operating software of the host device can either be android or IOS. The functionality developed is to record a ten-second sound, annotate it using a label filter from 1-5, as the number of the classes, and send this information alongside longitude, latitude, and DateTime to a database for saving.

As with all applications, a back-end was developed too, in order to handle the communication between the mobile application and the database. The backend is entirely developed in Python using Flask as a Webserver. Using Rest API technology, all the POST requests are received and handled by flattening and inserting the data into the database. For this purpose, a PostgreSQL database is also deployed.

Additionally, a visualization tool was also developed using the same back-end, enabling it to fetch the data from the database and visualize it in a map visual using google services. The visualization tool was created using React Native and Node JS too. The back-end again using Rest APIs technology, handles the GET requests and serves the data to the application by fetching them from the database.

4.1 Mobile Application – Frontend

The Frontend of the application is the user interface (UI) developed to allow users to send sounds. The main goal was to be fast and efficient and allow the users to send multiple sounds in a short period of time. For this purpose, the interaction between user and application was limited to the least button-pressing possible. The main interface is depicted in Figure 14, showing the UI as soon as the user initiates the application on their mobile phone.

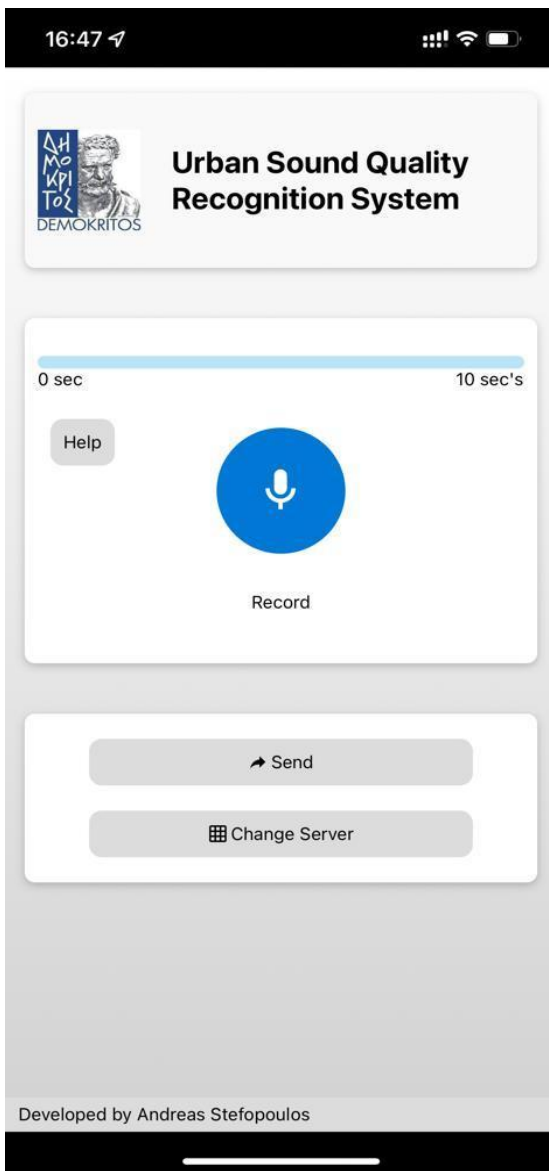


Figure 14: Main User Interface

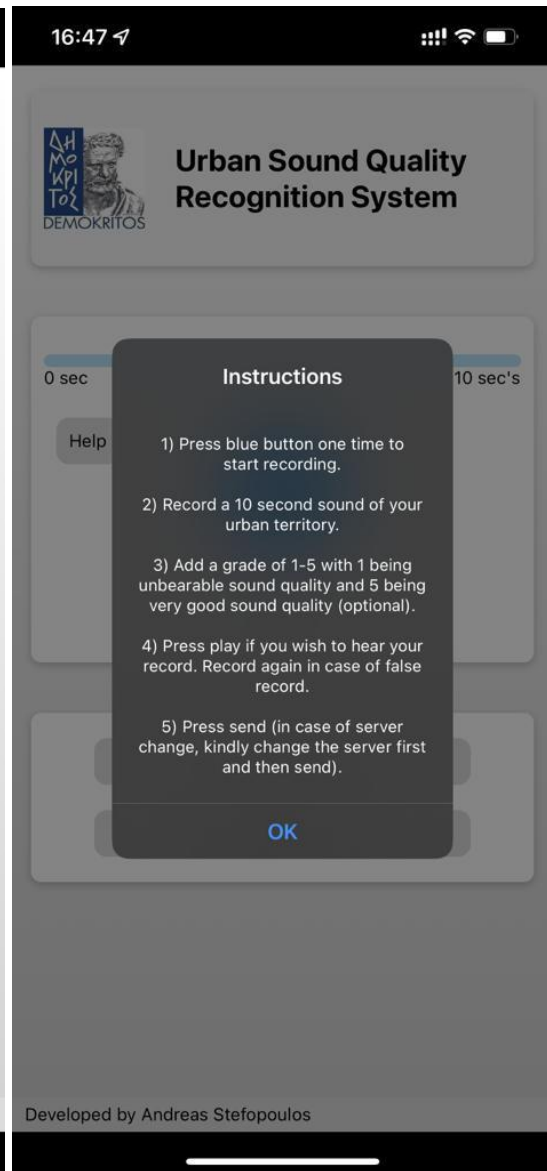


Figure 15: Help Button Instructions

The initiation of the recording can be achieved as soon as the user presses the blue button above “Record”. The process of recording starts and finishes automatically after ten seconds pass constructing a .wav file in the temporary memory of the phone. Before recording, the user may use the help button, which reveals a pop-up message showing the instructions for recording, and the logic followed is clarified in order to avoid mismatches in sounds and great deviations. Having a deviation between users is expected due to the different perceptions each one of us has.

The help message instructions are:

1. Press the blue button one time to start recording.
2. Record a 10 second sound of your urban territory.
3. Add a grade of 1-5 with 1 being unbearable sound quality and 5 being very good sound quality (optional).
4. Press play if you wish to hear your record. Record again in case of false record.
5. Press send (in case of server change, kindly change the server first and then send).

This message also reveals the functionality of the application. Continuing in the functionality, as soon as the user presses the record button, the UI changes, giving a new UI which can be seen in Figure 16. A bar is loading indicating the seconds elapsed while recording. There is a temporary functionality also here, converting the recording button to the stop button. While the stop button is pressed, the recording stops, and the whole record is deleted. This allows the users to initiate a new recording. After the elapse of ten seconds, the UI changes once more, revealing a new set of options with the label “Perceived Soundscape Quality” which is the way for the users to annotate the sound they have just recorded.

The selection for annotation can be done using a pop-up list appearing with labels “Optional” and numbers from 1 to 5. This pop-up list can be seen in Figure 19. As a next step for the user is to hit the send button and wait for the icon of send button to reappear (Figure 19). After each send, a pop-up message also appears, to notify the user that the recording has been sent successfully (Figure 20). In case of an issue, after pressing the send button (Figure 19), a pop-up message appears (Figure 21) to inform the user to contact the administrator. This mostly happens when the user is not connected to the internet or the server hosting the backend is not online.

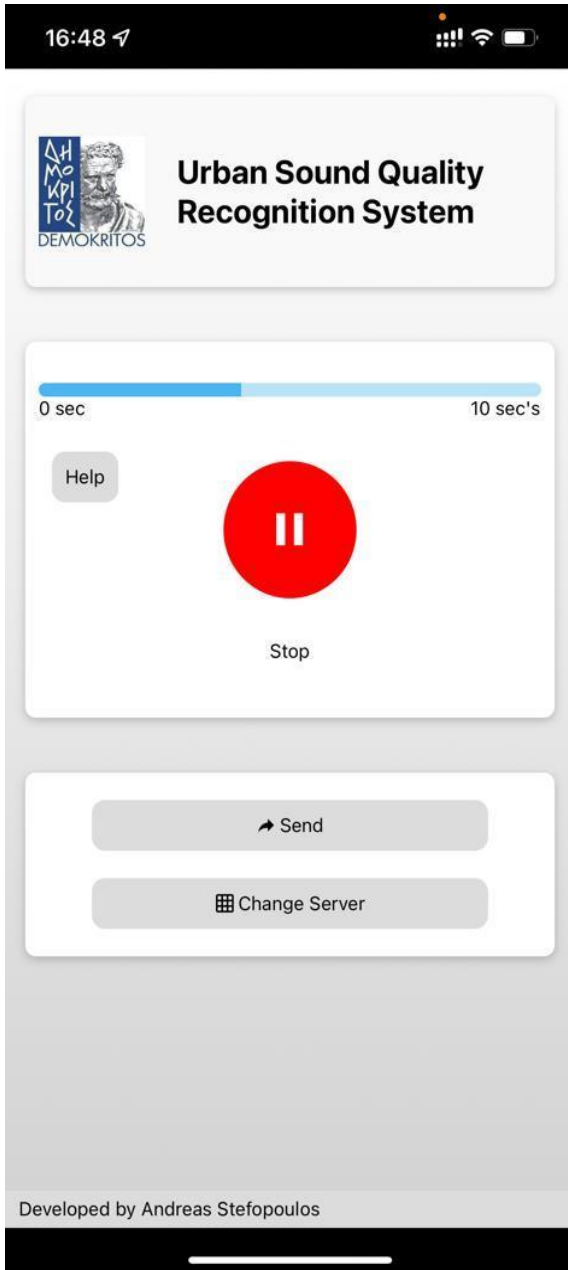


Figure 16: Recording in Progress

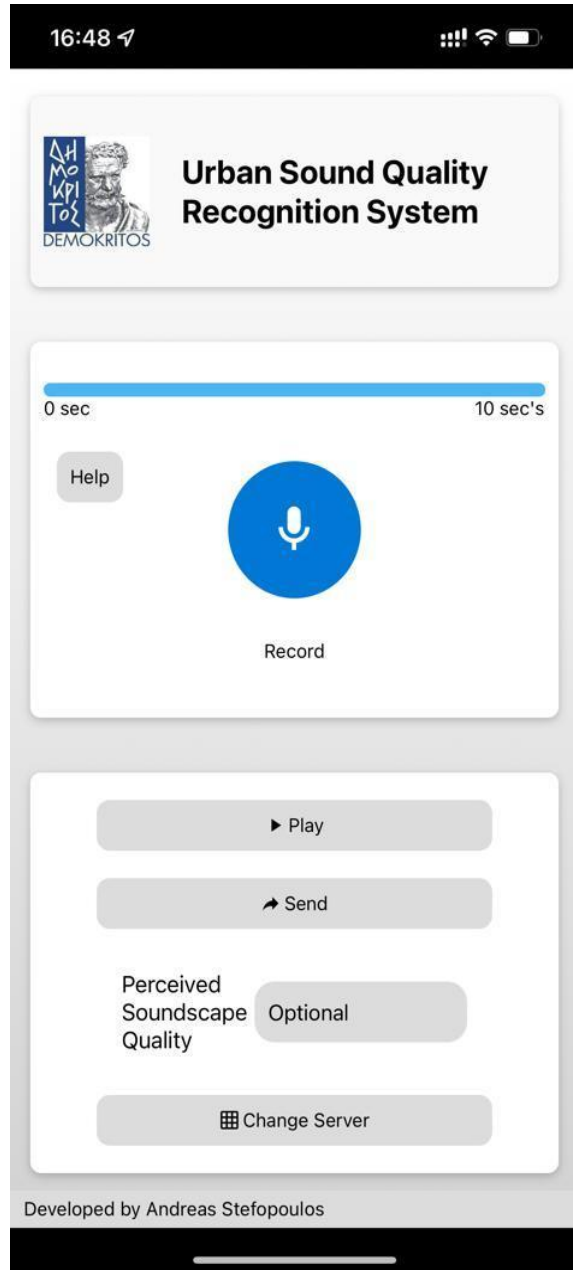


Figure 17: Recording Finished

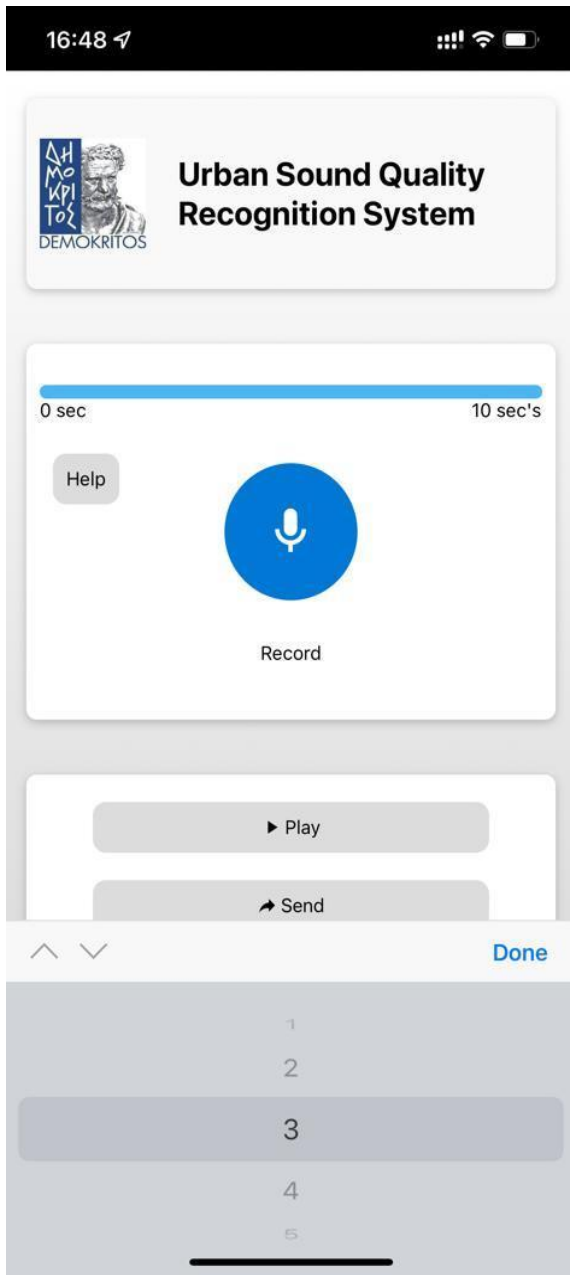


Figure 18: Selection List

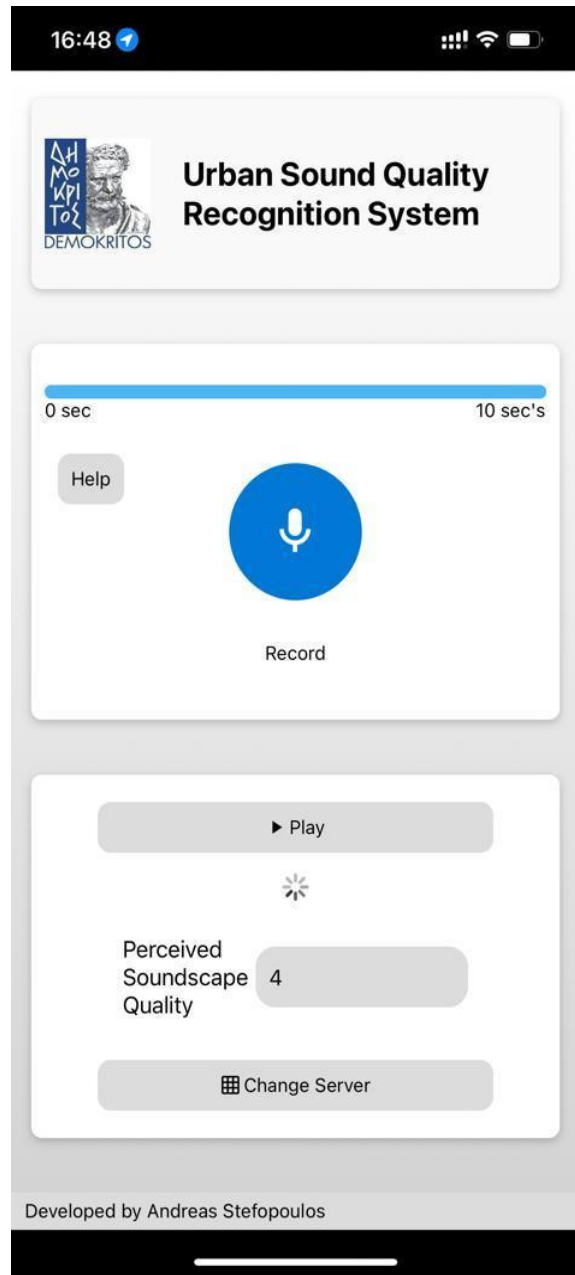


Figure 19: Sending the Record

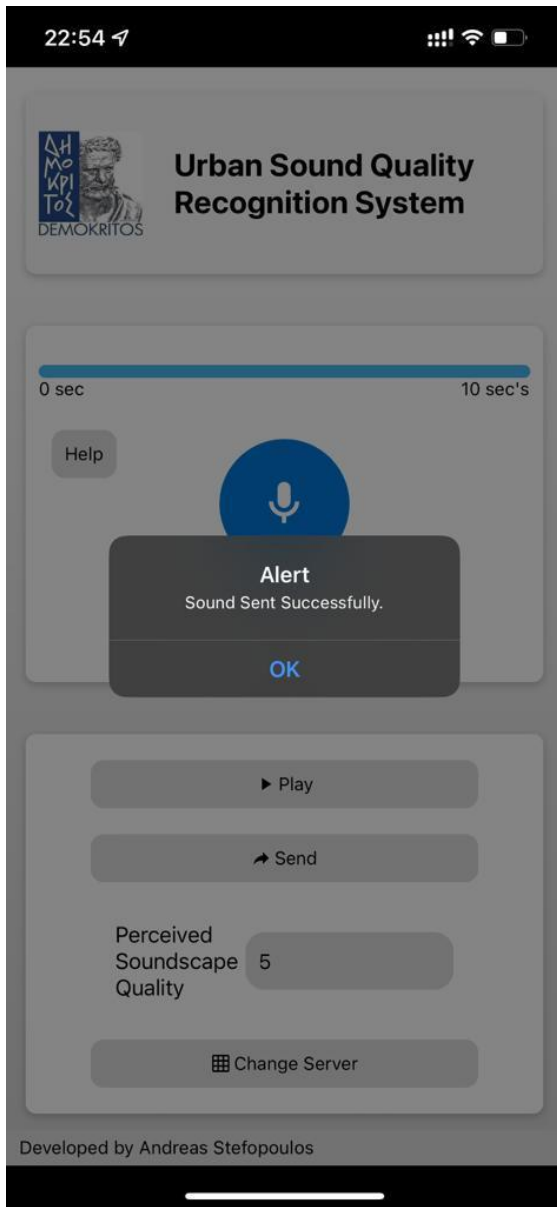


Figure 20: Recording Successfully Sent

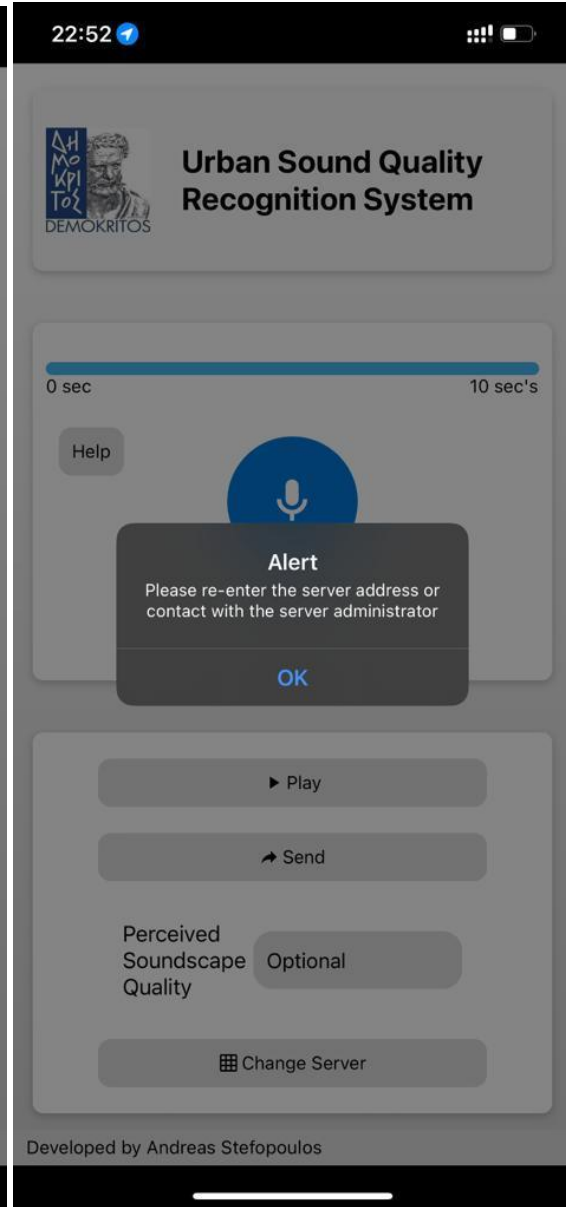


Figure 21: Recording not Sent Successfully

From the technical side, what happens inside the application is the construction of a 10-second .wav file stored temporarily. Additionally, when the user presses send, the location in longitude and latitude are also fetched alongside the DateTime timestamp. All this information is then formatted as JSON to be sent to the BackEnd of the application in the Webserver using REST APIs. Also, with the abovementioned JSON, the raw sound file in a .wav format is also sent to be handled from the WebServer. Lastly, the “change server” field is used in cases where the default server of the application is down, allowing

the user to modify the recipient of the API POST information, or in cases where the default server in the production server and the user want to send some recordings to a test server. This modification changes the server in the root code for the existing device and version of the application, being reset in case of reinstalling or updating.

4.2 Mobile Application – Backend

The mobile application sends the JSON file alongside the raw .wav file using APIs to the WebServer as mentioned in the previous chapter. The WebServer uses Flask as architecture and is deployed using Nginx and gunicorn, and accepts both GET and POST requests. The main processes done in the WebServer are two: the first is to handle the POST requests coming from the mobile application and the second is to handle the GET requests coming from the visualization part.

During the development of the WebServer, a database must be established and deployed as a migration to the process, in order for the WebServer to send the desired information directly to the database. For this purpose, a PostgreSQL database is used to receive the information and insert it into a table dedicated to the WebServer. For the migration between Flask and PostgreSQL was used the alembic library, which based on the arguments given, is produced the object to be included and called from the Flask in every insert of a new record.

After developing the part of the database and migrating to the WebServer, the process which will handle the contents of the requests was created using a decorator to set the URL which the WebServer will be available. In this process, the first action is to handle the raw file and save it in a static directory renaming the file as follows: “quality_deviceID_longitude_latitude_timestamp.wav”. The JSON data is flattened and converted to string values. These string values are then sent using the object created from the migration, which is called using the string values. When calling the database objects, an INSERT statement is taking place to insert the new record into the database. The columns of the table in the database have the following format:

1. index
2. name (Device ID)
3. sound_filename

4. data (the whole recording in base64 encoding for backup purposes)
5. latitude
6. longitude
7. sound_path (the path of the raw file where it was saved)
8. quality (annotation from the user)
9. predicted (field to be populated from the model which will be used)

The first process of handling the incoming data from the mobile applications ends with appending the new record of the abovementioned information to the database and storing the raw .wav file in a directory. For the purpose of visualization, we have the second decorator in the code, which leads to another URL for visualization purposes. The actions taken in this process are to fetch the data from the database and expose them using REST APIs. This results in responding to a GET request with all the data of the database coming from the website deployed. The website will be analyzed in the following section.

In order to deploy the WebServer, the domain “urbansoundsystem.ml” was purchased and re-routed using Route53 service from Amazon Web Services which was then included in sites-available in Nginx. For the deployment in production and to a physical server, instructions from Digital Ocean were followed.

Lastly, the whole infrastructure is developed in a way to have already embedded a client to call for deploying a PyTorch pre-trained model directly in the process.

4.3 Data Visualization

In the pipeline of the process of data gathering, visualization was also created, in order to monitor the progress of the sound recordings gathered and all the metadata stored in the database. For this purpose, a visualization was developed using React Native and Node JS to use the python backend to fetch and visualize data.

Having the back-end exposing the data existing in the database using APIs, the React Map is using GET requests to fetch the data from the database where data is exposed using an object created in a migration process.

After fetching the data, using Google API we are able to visualize the map of Google Maps and add pins for the location of each sound sent, based on the coordinates (Figure 22).

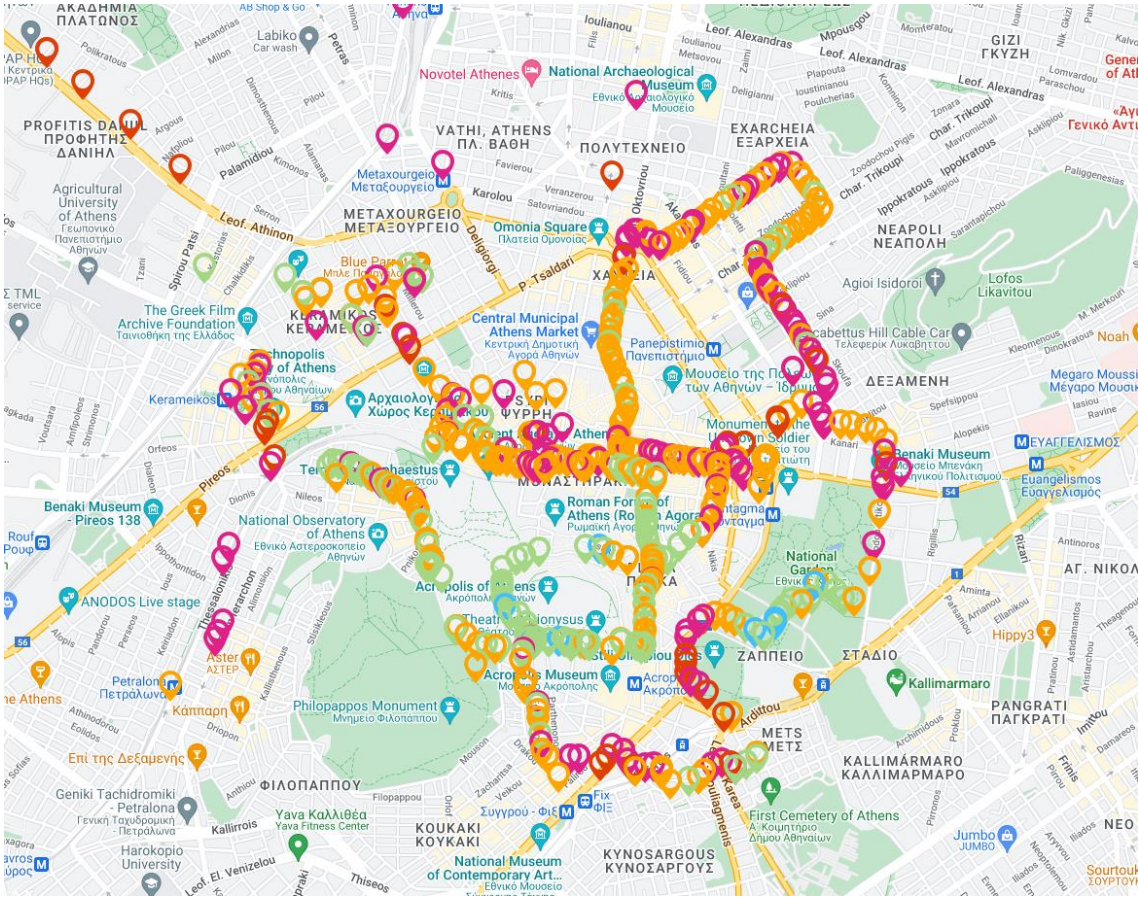


Figure 22: Example of Location Pins in Athens, Greece

Each location pin represents the quality annotated by the user. The interactive part of the visualization is a filter pane created in the top left corner and has the interface seen in Figure 23.

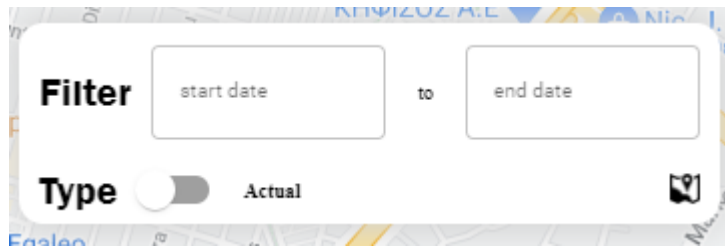


Figure 23: Filter Pane in React Map

By pressing the map icon in the bottom right corner of the filter pane, all information fetched from the database is shown, alongside all the information for color classification

and specific date filter applied information. By filtering January to August in the year 2022 in the pop-up map (Figure 24), followed by clicking the map icon we have the information shown in Figure 25.

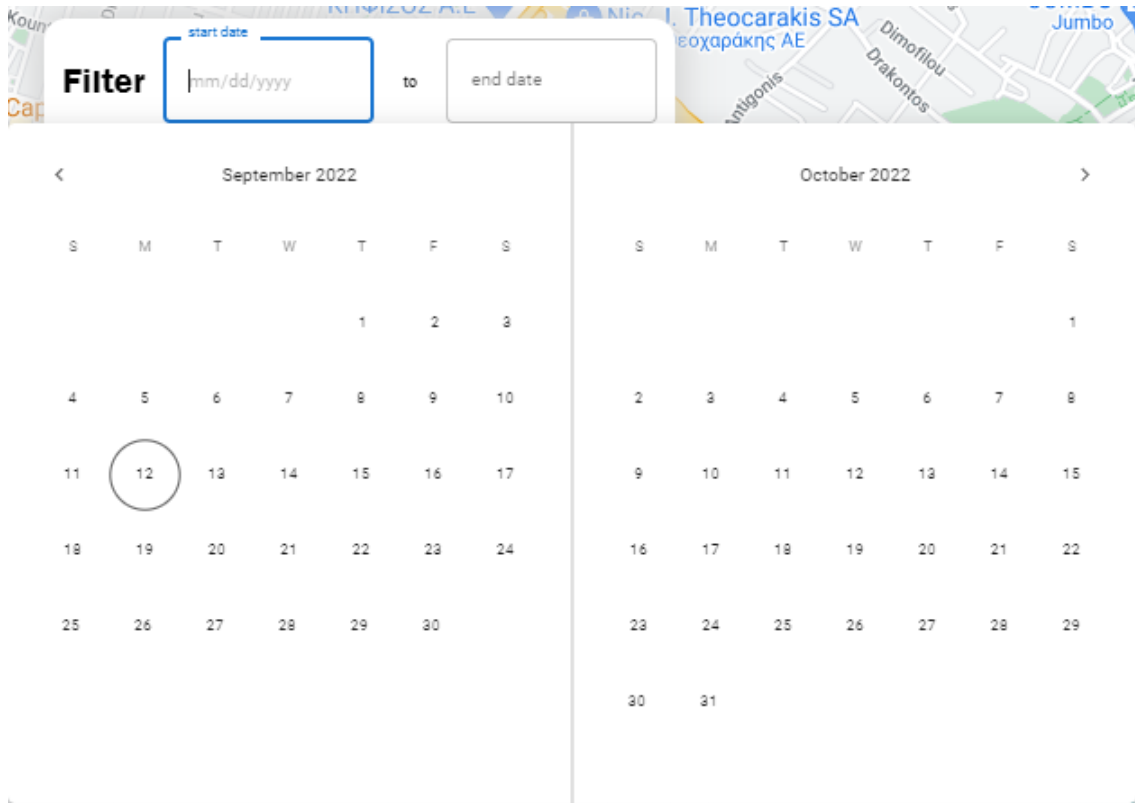


Figure 24: Date Filter Pop-Up

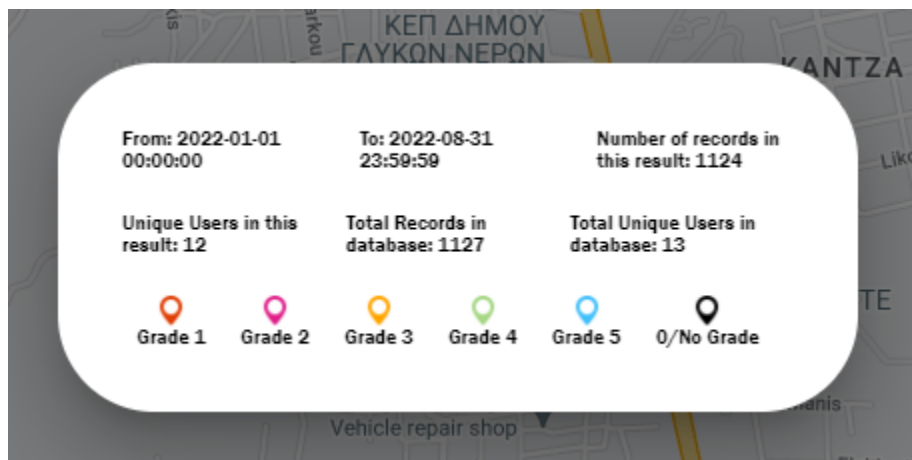


Figure 25: Details Button Result in Filtering

The metrics observed are from and to date period filtered, the number of records in the database for this result/filtering, the unique users (Device IDs) in this result/filtering, the numbers of total records in the database, and the total numbers of unique users in the database.

The color classification for the annotated value of every sound sent is represented with the icons in Figures 26-31.



Figure 26: Icon for Grade 0



Figure 27: Icon for Grade 1



Figure 28: Icon for Grade 2



Figure 29: Icon for Grade 3



Figure 30: Icon for Grade 4



Figure 31: Icon for Grade 5

Lastly, in the filter pane, there is also the option to switch between Actual and Predicted values, switching, in reality, to the column from the database we fetch the quality of the recordings. As mentioned before, there is the annotated value sent from the user and also another column “predicted” which will be populated from the PyTorch model.

4.4 Data Preparation - Backend

The data manipulation, after data has been gathered and stored in the remote server in the PostgreSQL database, is performed from an additional python script, created to fetch the data in a local directory, populating the directory with .wav files and creating the datasets for training and testing. This script also converts the data to mono channel (using pyffmpeg) and filters out any test files, comparing the name of the recordings with the records in the database.

Lastly, after every run of this script, a full backup is taking place, storing a .csv file with a database backup and all the raw files of the server stored locally in a specific directory.

5 Experiments

The experimentation of finding the optimal classification model is done in order to find optimal hyperparameters for the specific research and deploy the model in the end-to-end process, automating the training and prediction.

The experiments that took place are to train a model to classify the soundscape quality of a recording sent from users as accurately as possible. The main library used in these experiments is the `deep_audio_features` developed by Mr. Theodoros Giannakopoulos. The library uses PyTorch as the main deep learning library and the customization and hyperparameter tuning were implemented using the library's model as a baseline. The main Neural Network model used is a Convolutional Neural Network which after experimentations and tuning was changed to fit the data gathered.

As training inputs for the models were used the audio .wav recordings were gathered from the mobile application and prepared from the data backend script. The inputs for the models are mono-channel sound recordings in .wav format. The library's steps for the training process are to define the classes, `x_train`, `y_train`, `x_test`, and `y_test` and continue to generate a spectrogram for each sound and convert it to a Mel spectrogram using the `librosa` library.

In signal processing this approach is one of the most used, converting sound to an image representation, leading to a very, fast preprocess and also efficient feature extraction method. The preprocess implemented for feature extraction is the same in all the conducted experiments, generating for every sound file a Mel spectrogram with dimension 128x51 pixels.

5.1 Datasets

The data gathered from the developed backend were split and investigated using different approaches, to analyze better biases, optimizations, and best practices. The fact of rich in-context data gathered having all information regarding time, location, and device id help in analyzing the data from different perspectives.

The data was split into the following datasets:

- Dataset 1: A five-class dataset containing all data gathered from the mobile application, from all users, performing only a basic cleanup in test records. The data is split into train 80% and test 20%.
- Dataset 2: A five-class dataset consisting only of two IDs having validated that this data is correctly gathered and annotated. The first ID's records are located in Athens, Greece and the second ID's in Prague, Czech Republic. The data is split into two directories, country-oriented (cz, gr).
- Dataset 3: A three-class dataset containing all data gathered from the mobile application, from all users, performing only a basic cleanup in test records. Classes one and two are merged into class one, class three remained class three, four and five merged into class three. The data is split into train 80% and test 20%.
- Dataset 4: A three-class dataset consisting only of two IDs having validated that this data is correctly gathered and annotated. The first ID's records are located in Athens, Greece and the second ID's in Prague, Czech Republic. The data is split into two directories, country-oriented (cz, gr). Classes one and two are merged into class one, class three remained class three, four and five merged into class three.
- Dataset 5: The UrbanSound8K dataset was used for the transfer learning process. The original dataset consists of ten classes. In the current research, classes were merged to have a five-class dataset and a three-class dataset.

5.2 CNN Models

The main library used in the experiments is `deep_audio_features` created by Mr. Theodoros Giannakopoulos. Using this library, various PyTorch CNN models were created, to investigate the best-fitting ones.

After many experiments, three model architectures were used in the validation and experimentation process.

5.2.1 Models Architecture

First model architecture:

Category	Value
Activation Function	Leaky Relu
Batch Normalization	Yes
CNN Layers	4
First Linear Layer Activation Function	Leaky Relu
First Linear Layer Dimensions	Output Dimensions, 1024
First Linear Layer Dropout	0.75
Kernel Size	5x5 (2-Dimensions)
Linear Layers	3
Max Pool	2 Dimensions, Kernel Size = 2
Padding	2
Second Linear Layer Activation Function	Leaky Relu
Second Linear Layer Dimensions	1024, 256
Second Linear Layer Dropout	0.5
Stride	1
Third Linear Layer Activation Function	Leaky Relu
Third Linear Layer Dimensions	256, Output Dimensions

Second model architecture:

Category	Value
Activation Function	Leaky Relu
Batch Normalization	Yes
CNN Layers	2
First Linear Layer Activation Function	Leaky Relu
First Linear Layer Dimensions	Output Dimensions, 1024

First Linear Layer Dropout	0.75
Kernel Size	5x5 (2-Dimensions)
Linear Layers	2
Max Pool	2 Dimensions, Kernel Size = 2
Padding	2
Second Linear Layer Activation Function	Leaky Relu
Second Linear Layer Dimensions	1024, Output Dimensions
Second Linear Layer Dropout	0.5
Stride	1

Third model architecture:

Category	Value
Activation Function	Leaky Relu
Batch Normalization	Yes
CNN Layers	3
First Linear Layer Activation Function	Leaky Relu
First Linear Layer Dimensions	Output Dimensions, 2048
Kernel Size	5x5 (2-Dimensions)
Linear Layers	2
Max Pool	2 Dimensions, Kernel Size = 2
Padding	2
Second Linear Layer Activation Function	Leaky Relu
Second Linear Layer Dimensions	2048, Output Dimensions
Second Linear Layer Dropout	0.5
Stride	1

6 Results

The results are the outputs of the experimentation of the abovementioned models and datasets. The main metric used for evaluating the fit of the models is macro average f1 score in validation.

6.1 Results of Experimentation

6.1.1 Model 1

Experiment Number	Datasets	Validation Macro Average f1 Score	Prior Class Naïve Classifier
1	Dataset 1	0.35	0.31
2	Dataset 2- Train with Athens validation with Prague	0.19	0.19
3	Dataset 2- Train with Prague validation with Athens	0.25	0.21
4	Dataset 3	0.50	0.18
5	Dataset 4- Train with Athens validation with Prague	0.39	0.32
6	Dataset 4- Train with Prague validation with Athens	0.44	0.30

6.1.2 Model 2

Experiment Number	Datasets	Validation Macro Average f1 Score	Prior Class Naïve Classifier
7	Dataset 1	0.38	0.31
8	Dataset 2- Train with Athens validation with Prague	0.21	0.19
9	Dataset 2- Train with Prague validation with Athens	0.26	0.21
10	Dataset 3	0.51	0.18

11	Dataset 4- Train with Athens validation with Prague	0.40	0.32
12	Dataset 4- Train with Prague validation with Athens	0.38	0.30

6.1.3 Model 3

Experiment Number	Datasets	Validation Macro Average f1 Score	Prior Class Naïve Classifier
13	Dataset 1	0.33	0.31
14	Dataset 2- Train with Athens validation with Prague	0.20	0.19
15	Dataset 2- Train with Prague validation with Athens	0.22	0.21
16	Dataset 3	0.52	0.18
17	Dataset 4- Train with Athens validation with Prague	0.42	0.32
18	Dataset 4- Train with Prague validation with Athens	0.37	0.30

6.2 Transfer Learning

Due to low macro average f1 scores in the validation, in order to achieve better scores and due to the fact of a small data sample and human error in recording, transfer learning was implemented using dataset 5 (UrbanSound8K). For this series of experiments, model2 was selected.

The results per model after implementing transfer learning and then training only the last layer, are the following:

Experiment Number	Datasets	Validation Macro Average f1 Score	Prior Class Naïve Classifier
19	Dataset 1	0.42 – 0.47 (max)	0.31
20	Dataset 2- Train with Athens validation with Prague	0.19	0.19
21	Dataset 2- Train with Prague validation with Athens	0.24	0.21

22	Dataset 3	0.50	0.18
23	Dataset 4- Train with Athens validation with Prague	0.36	0.32
24	Dataset 4- Train with Prague validation with Athens	0.36	0.30

In experiment number 19, the maximum value of 47% was achieved by tuning the number of nodes in the linear layers.

6.3 Further Analysis

Diving deeper into the data, there were some additional logics implemented for training and validation. These logics are time-independent while user-independent and user-dependent training, in order to avoid using audio from the same timeframe. This leads to training with specific days/times and testing with the remaining days/times. For this series of experiments, the best-fitting model was selected (model2) and after implementing also transfer learning.

The results are the following:

Experiment Number	Datasets	Validation Macro Average f1 Score	Prior Class Naïve Classifier
25	Dataset 2 - Train with Athens validation with Prague - User Dependent	0.25	0.19
26	Dataset 2 - Train with Prague validation with Athens - User Dependent	0.28	0.21
27	Dataset 4 - Train with Athens validation with Prague - User Dependent	0.42	0.32
28	Dataset 4- Train with Prague validation with Athens - User Dependent	0.48	0.30
29	Dataset 1 - User Independent	0.40	0.31
30	Dataset 3 - User Independent	0.38	0.18

6.4 Results Interpretation

The results from the models after the initial training, using the audio samples gathered via the mobile application, appear to have low scores due to the fact that the dataset used is diverse, including a lot of bias, and relies purely on the annotation of the perception of different human beings. It is also observed that after applying transfer learning, the accuracy of the prediction rises, leading to the result that the volume of data gathered is not sufficient to train a neural network-based model to be accurate enough. Further investigation using more data and additional data preprocess methods, such as filtering out specific sounds or adjusting the microphone sensitivity, is needed, in order to identify more precisely the reasons for low scoring.

Lastly, considering the scenario of applying the above research in a real-world case by having no limitations in gathering samples from a distinct city, in order to use them as training inputs in a model and by having user independency, the maximum macro average f1 score that can be achieved in a 5-class case is 0.47. This result occurred by training model2 with transfer learning, meaning that the volume of data needed to be trained on is high. The architecture consists of 2 CNN layers and 2 linear layers, with a dropout of 0.75 in the first linear layer and 0.5 in the second. Additionally, the stride is set to 1 and the padding to 2. The number of nodes in the linear layers is 1024 per layer. In the 3-class case, the maximum macro average f1 score was achieved by model3 and is 0.52. The data trained was the data gathered from the mobile application and there was no better score achieved by applying transfer learning. The architecture consists of 3 CNN layers and 2 linear layers, with a dropout of 0.5 in the second linear layer. Additionally, the stride is set to 1 and the padding to 2. The number of nodes in the linear layers is 2048 per layer.

7 Conclusion

In the current dissertation, a mobile application was created in order to gather 10-second audio samples and, by using a custom back-end, store the data in a database/server aiming to be used as training inputs in a deep learning algorithm that would classify the audio quality, based on the annotated value which relies on the perception of the user. Furthermore, additional scripts were created, supporting the functionality of this process as a pipeline, allowing to automate and making it integration-available for additional uses.

A future implementation would be to apply the infrastructure created in static areas in an urban territory using small, energy-efficient devices (for example raspberry pies) and apply temporal analysis, tracking the changes of various factors such as the type of land use, the level of traffic, the presence of certain types of businesses or facilities, public safety, transportation planning, the construction rate and the variety of animal species presence.

The goal is to assess and manage the urban sound environment, as well as enhance the health and standard of living of city dwellers and spot any changes that can have a negative impact.

We can improve the quality of our environment by utilizing the AI ecosystem.

8 Bibliography

[1] "Soundscape, Urban Sound Quality and Quality of Life: A Review" by L.M. de Borst, R.L. de Waard, and F.B. de Vries, published in the Journal of the Acoustical Society of America in 2013. This article reviews the relationship between urban sound quality, soundscape, and quality of life.

[2] "A Methodology for the Assessment of Urban Soundscapes" by G.L. Harrison, P.J. Bell, and J.D. Simon, published in the journal Environmental Pollution in 2005. This paper presents a methodology for evaluating the soundscape of urban environments.

[3] "The Role of Soundscapes in Urban Planning and Design" by R.L. de Waard and L.M. de Borst, published in the journal Landscape and Urban Planning in 2016. This article discusses the importance of considering soundscapes in urban planning and design, and provides guidance on how to do so.

[4] "The Perception of Sound" by J.D. Moore, published by Cambridge University Press in 2003. This book covers the psychological aspects of sound perception, including topics such as pitch, timbre, and loudness.

[5] "Auditory Perception: An Analysis and Synthesis" by R.J. Dooling and D.B. Rye, published by MIT Press in 1991. This book provides a comprehensive overview of auditory perception, including topics such as pitch, timbre, loudness, and auditory masking. It also discusses the role of auditory perception in language and speech.

[6] "Acoustic Scene Classification Using Convolutional Neural Networks and Multiple Databases" by Sharath Adavanne, Archontis Politis, Toni Heittola, and Tuomas Virtanen (2017) - This paper presents a convolutional neural network (CNN) for acoustic scene

classification, which is trained on multiple databases and evaluated on a public dataset with state-of-the-art performance.

[7] "Environmental Sound Classification with Convolutional Neural Networks" by Justin Salamon, Christopher Jacoby, and Juan Pablo Bello (2015) - This paper proposes a deep convolutional neural network (CNN) for environmental sound classification, which uses a spectrogram as input and achieves state-of-the-art performance on various datasets.

[8] "Joint Optimization of Front-End Processing and Convolutional Neural Network for Robust Speech Recognition" by Yong Xu, Jun Du, Li-Rong Dai, and Chin-Hui Lee (2014) - This paper proposes a joint optimization method for front-end processing and convolutional neural network (CNN) to improve robustness of speech recognition in noisy environments.

[9] "Deep Convolutional Neural Networks and Data Augmentation for Environmental Sound Classification" by Huy Phan, Dinh Tuan Tran, Dang Lam Pham, and Thanh Son Dao (2016) - This paper presents a deep convolutional neural network (CNN) for environmental sound classification, which uses data augmentation techniques to improve performance on a public dataset.

[10] "Environmental Sound Classification Using Convolutional Neural Networks with Logarithmic Mel-Spectrogram" by Zhaoyang Zhang, Jiqing Han, and Wenwu Wang (2016) - This paper proposes a convolutional neural network (CNN) for environmental sound classification, which uses a logarithmic Mel-spectrogram as input and achieves state-of-the-art performance on various datasets.

[11] "End-to-End Environmental Sound Classification Using a 1D Convolutional Neural Network" by Justin Salamon and Juan Pablo Bello. This paper proposes an end-to-end approach to environmental sound classification using a 1D convolutional neural network, achieving competitive results on multiple sound datasets.

[12] "Automatic Environmental Sound Recognition: Performance versus Computational Cost of Convolutional Neural Networks" by Bartłomiej Stasiak and Sławomir T. Wierzchoń. This paper investigates the performance versus computational cost of different CNN architectures for environmental sound recognition, finding that the use of dilated convolutions can improve the trade-off between accuracy and computational cost.

[13] "Jointly Optimizing Feature Extraction and Classification for Environmental Sound Classification Using Convolutional Neural Networks" by Ke Chen, Meng Yu, and Anwen Hu. This paper proposes a novel method for jointly optimizing feature extraction and classification for environmental sound classification using CNNs, achieving state-of-the-art performance on several sound datasets.

[14] "Transfer Learning for Audio-Based Music Classification Using Convolutional Neural Networks" by Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter (2020) - This paper investigates transfer learning for audio-based music classification using convolutional neural networks (CNNs) and shows that transfer learning can be effective in reducing the amount of training data required and improving performance on new, unseen datasets.

[15] "A Transfer Learning Approach for Audio Classification Using Convolutional Neural Networks" by Akshay Verma, Aashish Tyagi, and Gaurav Gupta. (2018) - This paper presents a transfer learning approach for audio classification using CNNs, and shows that transfer learning can improve performance on the target task by leveraging knowledge learned from a source task.

[16] "PyTorch: An Imperative Style, High-Performance Deep Learning Library, " by A. Paszke et al. (2019) - This paper introduces PyTorch, a high-performance deep learning library that provides an imperative programming style, making it easy to use and efficient. PyTorch has been widely adopted by the research community and is used in a variety of applications, ranging from computer vision to natural language processing.

[17] Tzanetakis, G., & Cook, P. (2002). "Musical genre classification of audio signals" - This paper proposes a set of audio features for musical genre classification, including MFCCs, spectral centroid, and zero-crossing rate.

[18] Lidy, T., Rauber, A., & Pertusa, A. (2010). "Evaluation of feature extraction methods for large-scale content-based analysis of audio" - This paper evaluates different audio feature extraction methods for large-scale content-based audio analysis, including spectral features, MFCCs, and chroma features.

[19] Zhu, C., Li, X., & Li, J. (2016). "Audio Event Detection Using Convolutional Neural Networks with Layer-Wise Context Expansion and Pooling" - This paper proposes a convolutional neural network (CNN) architecture for audio event detection, which includes layer-wise context expansion and pooling to capture both local and global audio features.

[20] "Feature Learning Based Deep Supervised and Convolutional Neural Networks for Speech Emotion Recognition" by Xu et al. (2018) - This paper proposes a supervised deep feature learning method using convolutional neural networks (CNNs) and deep neural networks (DNNs) for speech emotion recognition, which achieved high accuracy on a public dataset.

[21] "Convolutional Recurrent Neural Networks for Bird Audio Detection" by Stowell et al. (2018) - This paper proposes a system for bird audio detection using a convolutional recurrent neural network (CRNN) that combines CNNs and recurrent neural networks (RNNs), achieving state-of-the-art performance on a public dataset.

[22] "Deep Learning for Audio Signal Processing" by Li et al. (2020) - This review paper provides an overview of recent advances in deep learning methods for audio signal processing, including feature extraction, speech recognition, music analysis, and sound source separation.

[23] "Convolutional Neural Networks for Music Classification and Retrieval" by Choi et al. (2017) - This paper explores the use of convolutional neural networks (CNNs) for music classification and retrieval tasks, using log-mel spectrograms as input features. The proposed system achieved state-of-the-art performance on multiple music datasets.

[24] "Mel-spectrogram and inverse mel-spectrogram for speech enhancement" by Xu et al. (2014) - This paper proposes the use of mel-spectrograms for speech enhancement, a technique for improving the quality of speech in noisy environments. The authors found that their proposed method achieved state-of-the-art performance on multiple speech enhancement tasks.

[25] "Environmental sound recognition using sparse representations of spectrograms" by Yang et al. (2017) - This paper proposes a sparse representation approach to environmental sound recognition, using spectrograms as input features. The authors found that their proposed method achieved state-of-the-art performance on multiple environmental sound datasets.

[26] "Deep learning for music genre classification: A comparative analysis" by P. Khorrami et al. (2019) - This paper compares the performance of several deep learning models for music genre classification, and finds that a combination of convolutional and recurrent neural networks with log-mel spectrograms as input achieves the best results.

[27] "Convolutional neural networks for music classification: a comparison study" by L. F. Oliveira et al. (2018) - This paper compares the performance of several convolutional neural network architectures for music classification, and concludes that a deep CNN with log-mel spectrograms as input achieves state-of-the-art results.

[28] "A Deep Learning Approach for Speech-to-Song Conversion Using Log-Mel Spectrograms" by M. Yoon et al. (2019) - This paper proposes a deep learning approach

for converting speech to singing voice using log-mel spectrograms, and achieves high-quality results on various datasets.

[29] "A systematic study of the class imbalance problem in convolutional neural networks" by G. Garcia et al. (2019) - This paper investigates the effect of class imbalance on the performance of CNNs and provides recommendations for hyperparameter tuning to improve their performance on imbalanced datasets.

[30] "Learning features of music from scratch" by Jansson et al. (2017) - This paper presents a method for learning features from raw audio data, which involves converting the audio signals to Mel spectrograms and using a convolutional neural network (CNN) to learn the features. The proposed method achieved state-of-the-art performance on several music classification tasks.

[31] "A comparative study of CNN-based deep learning architectures for recognizing environmental sounds" by Sengupta and Biswas (2018) - This paper compares different CNN-based deep learning architectures for environmental sound recognition using mel-spectrograms as input and evaluates their performance on a standard dataset.

[32] "Audio-based Bird Species Identification using Deep Convolutional Neural Networks" by Stowell et al. (2018) - This paper proposes a deep convolutional neural network (CNN) for bird species identification using mel-spectrograms as input and achieves state-of-the-art performance on a challenging dataset.

[33] "Attention-based Convolutional Neural Networks for Acoustic Event Detection" by Xu et al. (2017) - This paper proposes an attention-based CNN for acoustic event detection using mel-spectrograms as input, which achieves state-of-the-art performance on a standard dataset.

[34] "Event Detection and Classification using Convolutional Neural Networks" by A. Mesaros et al. (2016) - This paper proposes a CNN-based system for audio event detection and classification, which uses log-mel spectrograms as input and outperforms state-of-the-art approaches on several datasets.

[35] "Multi-label classification of acoustic scenes using convolutional neural networks" by D. Stowell et al. (2015) - This paper proposes a CNN-based system for multi-label classification of acoustic scenes, which uses log-mel spectrograms and achieves state-of-the-art performance on a dataset with 15 acoustic scene categories.

[36] "A deep convolutional neural network approach for environmental sound classification" by Y. Xu et al. (2016) - This paper presents a CNN-based approach for environmental sound classification, which uses log-mel spectrograms and achieves state-of-the-art performance on a large-scale dataset with 50 categories.

[37] "CNN-Based Acoustic Event Detection with Multi-Scale Multi-Task Learning" by X. Xu et al. (2020) - This paper proposes a multi-scale multi-task CNN-based system for acoustic event detection, which uses log-mel spectrograms and achieves state-of-the-art performance on a dataset with 10 acoustic event classes.

[38] "Attention-based convolutional neural networks for acoustic scene classification" by Phan et al. (2017) - This paper proposes an attention-based CNN model for acoustic scene classification, which can selectively focus on different parts of the input spectrogram to improve classification performance.

[39] "A transfer learning approach for environmental sound classification using small datasets" by Kuleshov et al. (2017) - This paper proposes a transfer learning approach for environmental sound classification, which uses a pre-trained CNN model and fine-tunes it on small datasets to improve classification performance.

[40] "Audio Event Recognition Using Deep Learning with Time-Frequency Representations" by Y. Kong and J. Smith (2018) - This paper presents a CNN-based approach for audio event recognition, which uses log-mel spectrograms as input and incorporates a time-frequency feature aggregation module to capture both temporal and spectral information.

[41] "Convolutional Recurrent Neural Networks for Polyphonic Sound Event Detection," by Y. Xu et al. (2017) - This paper proposes a deep CNN-RNN architecture for polyphonic sound event detection, which uses log-mel spectrograms as input and incorporates both convolutional and recurrent layers to model the temporal and spectral dependencies in the audio signals.

[42] "Music Genre Classification using Convolutional Neural Networks," by S. Zhang et al. (2017) - This paper presents a CNN-based approach for music genre classification, which uses log-mel spectrograms as input and incorporates both convolutional and pooling layers to extract features from the audio signals.

[43] "A Multitask Learning Approach for Audio Event Classification Using Convolutional Neural Networks," by T. Heittola et al. (2018) - This paper proposes a multitask learning approach for audio event classification, which uses a deep CNN to extract features from log-mel spectrograms and performs both event classification and tagging simultaneously.