



Protecting Copyright Ownership via Identification of Remastered Music in Radio Broadcasts

M. A. P. P. Marasinghe

Index No : 15000877

Supervisor : Dr. K.L. Jayaratne

Co-supervisor : Dr. M.I.E. Wickramasinghe

February 2020

Submitted in partial fulfillment of the requirements of the
B.Sc in Computer Science Final Year Project (SCS4124)



Protecting Copyright Ownership via Identification of Remastered Music in Radio Broadcasts

M.A.P.P. Marasinghe

Declaration

I certify that this dissertation does not incorporate, without acknowledgement, any material previously submitted for a degree or diploma in any university and to the best of my knowledge and belief, it does not contain any material previously published or written by another person or myself except where due reference is made in the text. I also hereby give consent for my dissertation, if accepted, be made available for photocopying and for interlibrary loans, and for the title and abstract to be made available to outside organizations.

Candidate Name : M.A.P.P. Marasinghe

.....
Signature of Candidate

Date :

This is to certify that this dissertation is based on the work of

Mr. M.A.P.P. Marasinghe

under my supervision. The thesis has been prepared according to the format stipulated and is of acceptable standard.

Supervisor Name : Dr. K.L. Jayaratne

.....
Signature of Supervisor

Date :

Co-supervisor Name : Dr. M.I.E. Wickramasinghe

.....
Signature of Co-supervisor

Date :

Abstract

Preface

Acknowledgement

Table of Contents

Declaration	i
Abstract	ii
Preface	iii
Acknowledgement	iv
List of Figures	vii
List of Tables	viii
List of Acronyms	ix
1 Introduction	1
1.1 Background to the Research	1
1.2 Research Problem and Research Questions	4
1.2.1 Project Aim	4
1.2.2 Research Questions	4
1.2.3 Objectives	5
1.3 Methodology	5
1.4 Outline of the Dissertation	9
1.5 Scope and Delimitations	9
1.5.1 In Scope	9
1.5.2 Out Scope	9
1.6 Conclusion	10
2 Literature Review	11
3 Design	12
3.1 Architectural Design	12
3.2 Principle Component Analysis (PCA)	13
3.2.1 PCA with Raw Dataset	13
3.2.2 PCA with Dataset Normalized by Z-score	14
3.2.3 PCA with Dataset Normalized by Rescaling	14
3.3 Scale Invariant Feature Transform (SIFT) Based Approach	14

4 Implementation	15
5 Results and Evaluation	16
5.1 Experiments	16
6 Conclusions	17
References	18

List of Figures

1.1	Key controlling parameters of STFT[1]	1
1.2	Architecture of the existing system	2
1.3	Extracting peaks and generating a hash value[1]	3
1.4	Remastered Song Identification Process	5
1.5	Key parameters on Short Time Fourier Transformation (STFT). 2048 bits long window with 1024 bits long overlapping area.	6
1.6	Generated colour image of spectrogram after preprocessing.	6
1.7	Spectrogram transformations on audio enhancements. (a) is the spectrogram image of a original song. (b) 20% pitch increase, (c) 20% pitch decrease, (d) 20% tempo increase and (e) 20% tempo decrease spectrogram images.	7
1.8	Average Accuracy Values for Different Threshold Values	9
3.1	Architectural Design	12
3.2	PCA coefficients weighted by eigen values	13
3.3	PCA coefficients weighted by eigen values (Normalized by Zscore) .	14
3.4	PCA coefficients weighted by eigen values (Normalized by Rescaling)	14

List of Tables

List of Acronyms

BLOB Binary Large Object

DoG Difference of Gaussians

OSCA Outstanding Song Creators Association

PCA Principle Component Analysis

SIFT Scale Invariant Feature Transform

STFT Short Time Fourier Transformation

Chapter 1

Introduction

1.1 Background to the Research

According to the intellectual property act of Sri Lanka[2], royalties must be paid to the original artistes when a song is broadcast on a radio channel. Each radio channel is maintaining a playlist to keep track of the songs that were broadcast throughout the day. That playlist can later be used to pay royalties to the respective artistes. However, in order to streamline and regulate the royalty payment process, it is vital to have a method to monitor the radio broadcasts. Manual radio broadcast monitoring is infeasible and expensive due to increasing number of both radio channels and songs. In manual monitoring a person should be assigned to each channel who needs to keep record of each song in the radio broadcast of that assigned channel. Due to the increasing number of songs and the fallible nature of humans such a monitoring task is prone to errors and inaccuracies. Hence an automated radio broadcast monitoring approach must be considered as an viable alternative in the modern day radio broadcast monitoring.

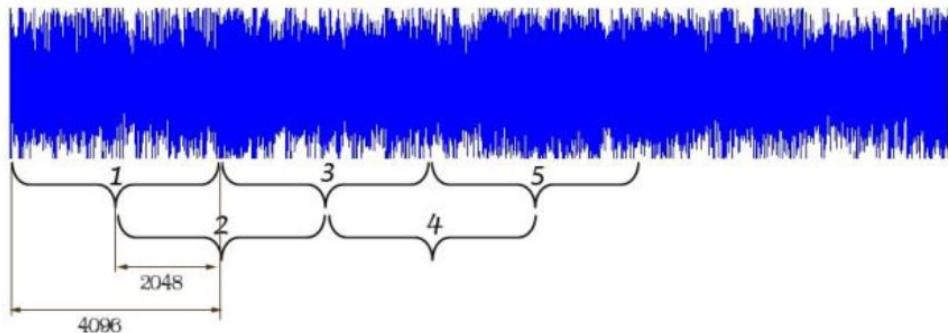


Figure 1.1: Key controlling parameters of STFT[1]

In the research “Radio Broadcast Monitoring to Ensure Copyright

Ownership”[1], researchers have implemented an automated radio broadcast monitoring system (refer the Figure 1.2 for the architecture) which has achieved 97.14% overall accuracy in identifying original songs in radio broadcasts. The researchers introduced an audio fingerprint to register and identify songs. The fingerprint was introduced as a series of hash values extracted from frequency domain audio signal. Time domain signal was converted to frequency domain by using STFT, which used 4096 bits long window and 2048 bits long overlapping area as shown in Figure 1.1. Then five peak values were extracted for each window by dividing mid frequency level into five bins and taking peak value from each bin. Extracted five peak values were used to create a hash value as depicted in Figure 1.3.

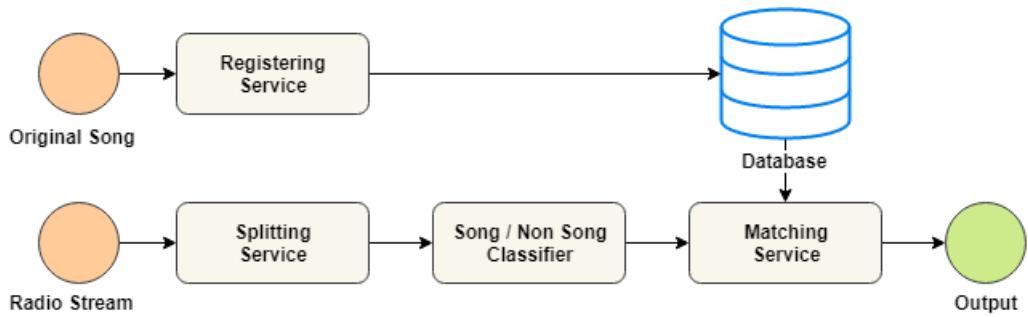


Figure 1.2: Architecture of the existing system

In contemporary radio broadcasts, channels tends to alter songs by including commercials and dialogues and by remastering the original song. Remastering can be done by adding or subtracting elements, or by changing pitch, equalization, dynamics or tempo[3]. Even though the above mentioned radio broadcast monitoring system’s accuracy is not significantly affected by commercials and dialogues included in songs, the system is unable to identify a song when that song is remastered by the radio channel as changing pitch, equalization, dynamics or tempo which directly affects both time domain and frequency domain audio data.

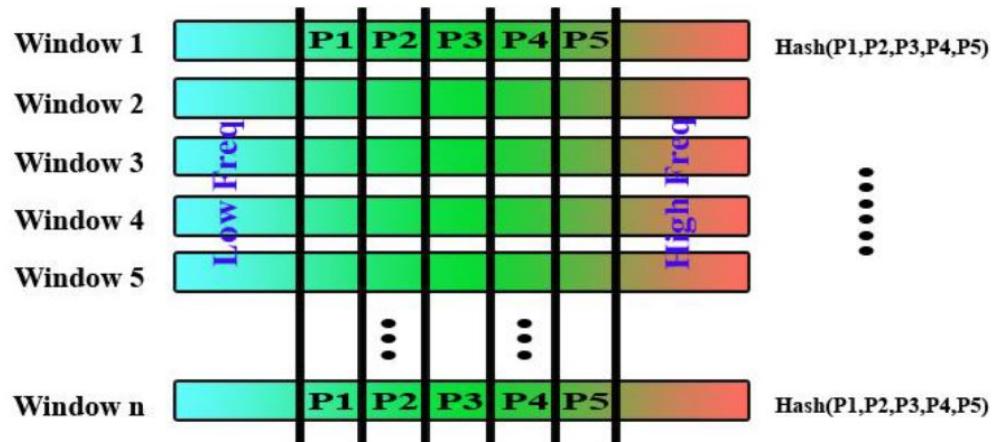


Figure 1.3: Extracting peaks and generating a hash value[1]

Timbre, tempo, timing, structure, key, harmonization and lyrics are the basic musical facets that can be identified[3]. Timbre, also known as tone colour is the music facet which makes a difference of different sound productions even when they have the same pitch and loudness. Simply it is what makes a difference between a piano and a violin playing the same note at the same volume. Timbre can be changed due to the use of different sound enhancing and processing techniques or to the use of different instruments and configurations. Tempo is the speed or pace of the music which can be easily changed by playing the music in different speeds. The music facet of timing is rhythmic structure of the music which can be altered by the changes to the drum section. Structure is the arrangement of music sections, and music structure alterations can be made while remastering. Key, harmonization and lyrics are tonality, chords and words of the music which can be altered while remastering.

In order to identify remastered music in radio broadcasts, existing literature on cover song identification and music similarity measures can be used as foundation study to this research. Directly implementing a cover song identification method or a music similarity measure to identify remastered music in radio broadcasts is not possible as there is limited time to do the identification and it is not just comparing two music clips to find similarity, but comparing a radio broadcast with more than twenty thousand song database.

1.2 Research Problem and Research Questions

1.2.1 Project Aim

Aim of this research is to utilize computational theories and tools to protect copyright ownership of artistes in radio broadcasts.

1.2.2 Research Questions

Main three research questions are identified to address the challenges in music identification when remastered songs are broadcasted in radio.

1. What are forms of alterations to the basic musical facets in remastered music?
2. What are the approaches of identifying remastered music?
3. What approach can be used to identify remastered music in radio broadcasts?

1.2.3 Objectives

Answers to above research questions are obtained by accomplishing the five objectives.

1. Gather and identify the alterations made to remastered music when compared with the original music.
2. Review existing cover song identification methods and music similarity measures to implement feature extraction methods for relevant features.
3. Identify similarity descriptive features with respect to the identified forms of alterations.
4. Introduce a new music similarity descriptor using identified features.
5. Use introduced music similarity descriptor to identify remastered music in radio broadcasts.

1.3 Methodology

In the proposed method of remastered song identification, various algorithms are used to extract the audio features, create audio descriptors and match against stored descriptors. Hence we have divided our remastered song identification process in to five steps.

1. Preprocessing
2. Feature Extracting
3. Descriptor Storing (Registering)
4. Matching
5. Postprocessing

Processes of the above steps will be discussed in the following subsections.

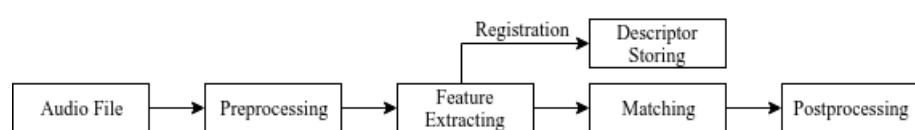


Figure 1.4: Remastered Song Identification Process

Preprocessing

In default audio data is represented in the time domain. Since even a small change in an audio changes the time domain representation drastically, using the time domain representation of the audio to extract features is not recommended. Hence time domain audio signal is converted to frequency domain signal by using STFT method. STFT is a sequence of Fourier transforms of a windowed signal[4]. 2048 bits long window with 50% overlapping was used as STFT key parameters as depicted in Figure 1.5.

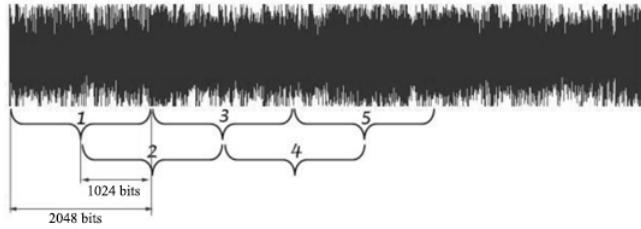


Figure 1.5: Key parameters on STFT. 2048 bits long window with 1024 bits long overlapping area.

STFT is often visualized using its spectrogram[4], which is an intensity plot of STFT magnitude over time. The generated spectrogram is converted to a color image as shown in Figure 1.6. Axis labels and ticks are removed to stop identification of them as key points in feature extracting step.

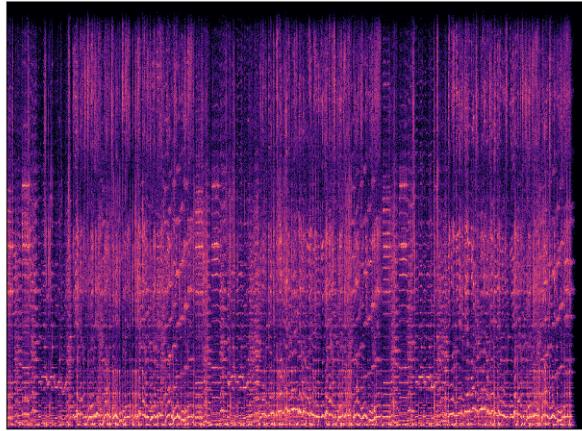


Figure 1.6: Generated colour image of spectrogram after preprocessing.

Feature Extracting

STFT spectrogram itself can be considered as an audio descriptor[5]. This method uses Scale Invariant Feature Transform (SIFT)[6] to extract the features which are robust to music remastering. In the Figure 1.7, it can be observed that when tempo

is altered the spectrogram will either expand or compress with the time axis and when pitch is altered the spectrogram will either shift upwards or downwards with the frequency axis.

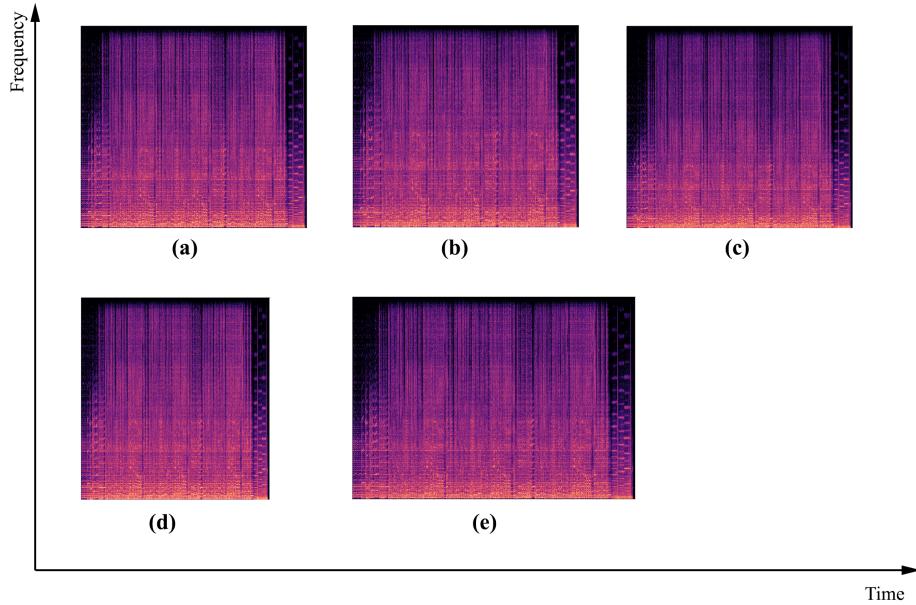


Figure 1.7: Spectrogram transformations on audio enhancements. (a) is the spectrogram image of a original song. (b) 20% pitch increase, (c) 20% pitch decrease, (d) 20% tempo increase and (e) 20% tempo decrease spectrogram images.

SIFT is used in computer vision to identify scale invariant features of an image. SIFT features are invariant to image rotation, scale alterations and illumination[6]. The SIFT feature extractor used in this method consists of four main steps.

1. Scale space extrema detection: Gaussian filters of different scales are applied to the image and potential key points are selected as local minima or maxima of the Difference of Gaussians (DoG) for multiple scales.
2. Keypoint localization: Keypoints that have low contrast or those that are poorly located along edges are filtered out.
3. Orientation assignment: One or more orientations are assigned to each keypoint based on local image gradient.
4. Keypoint descriptor generation: Orientation histograms are created for 4×4 pixel neighborhoods for each keypoint. Each histogram consists 8 bins, hence $(4 \times 4 \times 8)$ 128 dimensional descriptor is generated.

A Set of extracted 128 dimensional descriptors works together in describing the input audio file. Extracted SIFT features are invariant to image stretch and translation which makes them better features to be used in audio identification algorithm which is robust to tempo alterations and pitch shifting.

Descriptor Storing (Registering)

SIFT descriptors of original songs must be stored to use them in the matching step of the remastered song identification process. Generally 3-5 minute music clip will have around 2000 key points in its STFT spectrogram. Hence a 2000×128 matrix will be generated for each original song that will be registered.

The descriptor matrix of each original song is converted to a binary string and that binary string is stored in the database as Binary Large Object (BLOB)s. Converting to binary string and storing the matrix as a BLOB will ensure fast recreation of the matrix while retrieving[7].

Matching

Flann KD Tree Matching [8].

Postprocessing

the most similar song and matched keypoint count for a given query audio clip is identified in the matching step. But it doesn't exactly mean that query audio clip contains that song. Because the number of keypoints that were matched represents how much the query song matched to the most similar song. Hence there should be a threshold keypoint count to determine whether a query audio clip contains a song in our database or not. But using just a threshold value won't work here since different query audio clips generate different number of key points to match against the database. Hence ratio based threshold is recommended as a measure to determine whether the matched song is actually a correct match. Keypoint ratio can be obtained by the below equation.

$$\text{Keypoint ratio} = \frac{\text{Matched keypoint count}}{\text{Keypoints generated for query audio clip}}$$

Based on this keypoint ratio, a threshold is used to determine the validity of the match found. In order to find this threshold value we have used 844 different audio clips with variable durations to match against 2300 original songs. Those 844 audio clips had 519 audio clips which had songs and 325 audio clips which didn't have songs from those 2300 original songs. And we calculated accuracy for 18 testcases which will be discussed in section 5.1, and took average of those 18 accuracies for variable threshold values. Then results were illustrated as shown in the Figure 1.8.

Global peak can be observed in the illustration which makes that value a clear

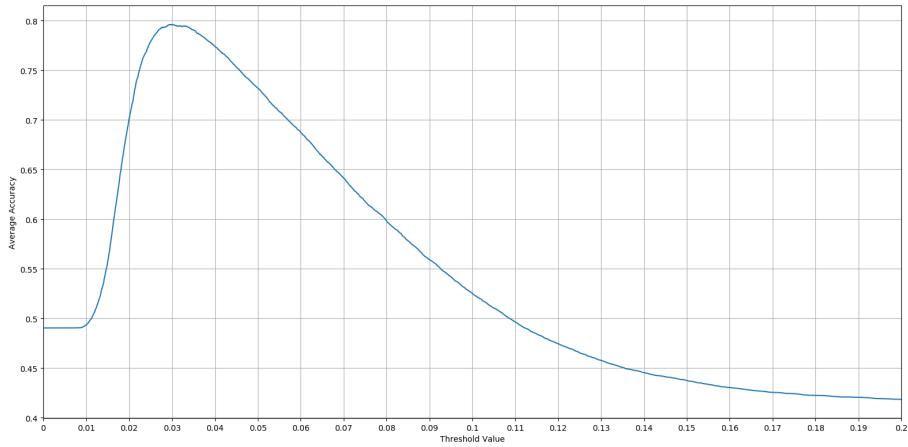


Figure 1.8: Average Accuracy Values for Different Threshold Values

threshold point. 0.0298 is the threshold value that was found. Hence if keypoint ratio of a query audio is larger than 0.0298 then it's identified as a valid match to the song that was identified in the matching step, otherwise it's identified as a invalid match. This threshold point makes this method to clearly identify whether a query audio has a song which is in a database or not.

1.4 Outline of the Dissertation

1.5 Scope and Delimitations

1.5.1 In Scope

The following areas will be covered under the research project.

- Exploration of possible remastering techniques and outcomes.
- Introduction of music similarity descriptor.
- Identification of remastered music in radio broadcasts.

1.5.2 Out Scope

The following areas will not be covered under the research project.

- Identification of instrumental, acoustic or medley covers of original music.
- Identification of quotations in music such as lyrical quotations or musical quotations.

1.6 Conclusion

Chapter 2

Literature Review

Chapter 3

Design

3.1 Architectural Design

Basic framework to identify music in radio broadcasts is proposed in “Radio Broadcast Monitoring to Ensure Copyright Ownership”[1]. And it’s currently implemented and deployed in Outstanding Song Creators Association (OSCA). Hence the infrastructure to contain the basic framework of a radio monitoring system is already there. When different methods of music identification are applied to radio broadcast monitoring, only registering service, database and matching service will be changed conserving the other modules that are in the architecture as shown in the Figure 3.1.

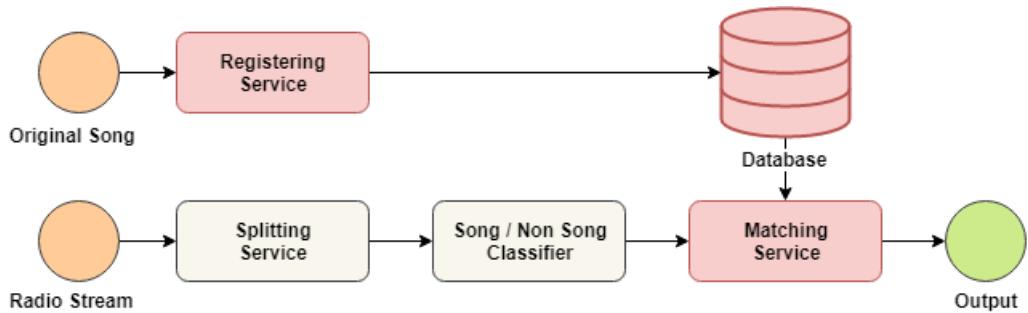


Figure 3.1: Architectural Design

Registering service takes original songs as input and generate specific descriptors to be stored in the database. Database stores the generated audio descriptors by registering service in retrieval friendly framework to help matching service to retrieve descriptors faster. Matching service takes a query audio clip and determine whether that query audio clip contains song registered before.

3.2 Principle Component Analysis (PCA)

There are many different approaches to identify music by extracting different audio features as discussed in the Chapter 2. Since there are very low number of researches conducted on sinhala music identification, finding audio features which can differentiate two sinhala songs was required to continue the research. Hence Principle Component Analysis (PCA) was conducted to find features on a 5000 song dataset extracting 27 different audio features and results were collected for different normalization techniques.

3.2.1 PCA with Raw Dataset

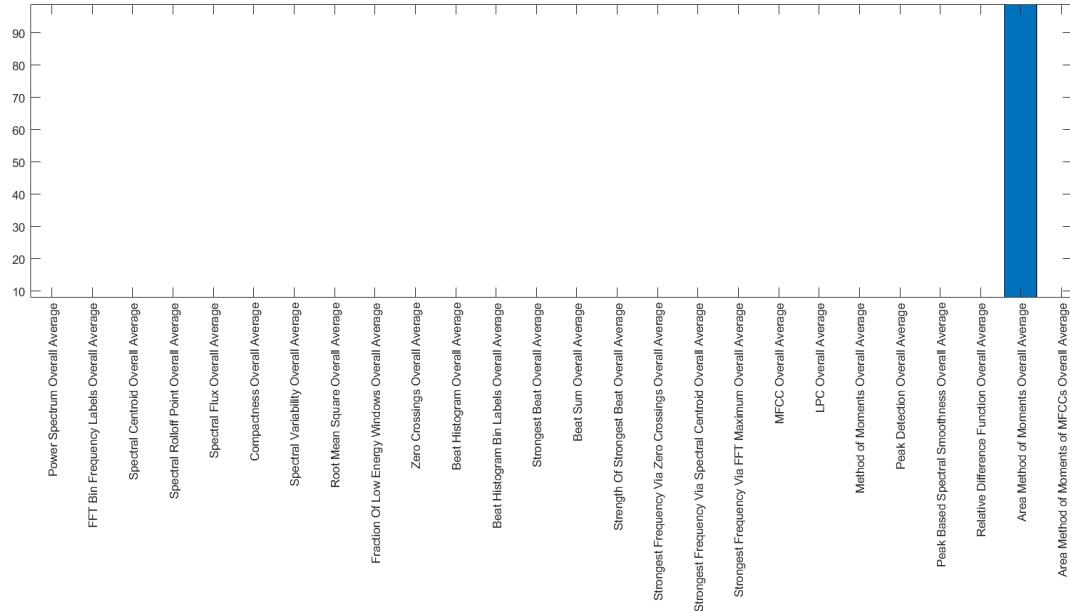


Figure 3.2: PCA coefficients weighted by eigen values

3.2.2 PCA with Dataset Normalized by Z-score

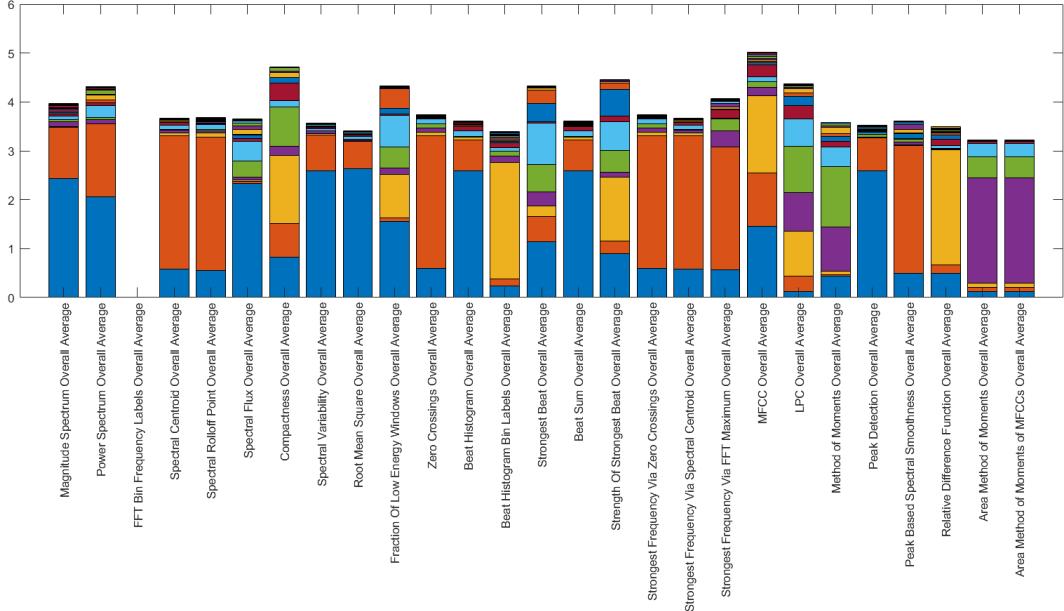


Figure 3.3: PCA coefficients weighted by eigen values (Normalized by Zscore)

3.2.3 PCA with Dataset Normalized by Rescaling

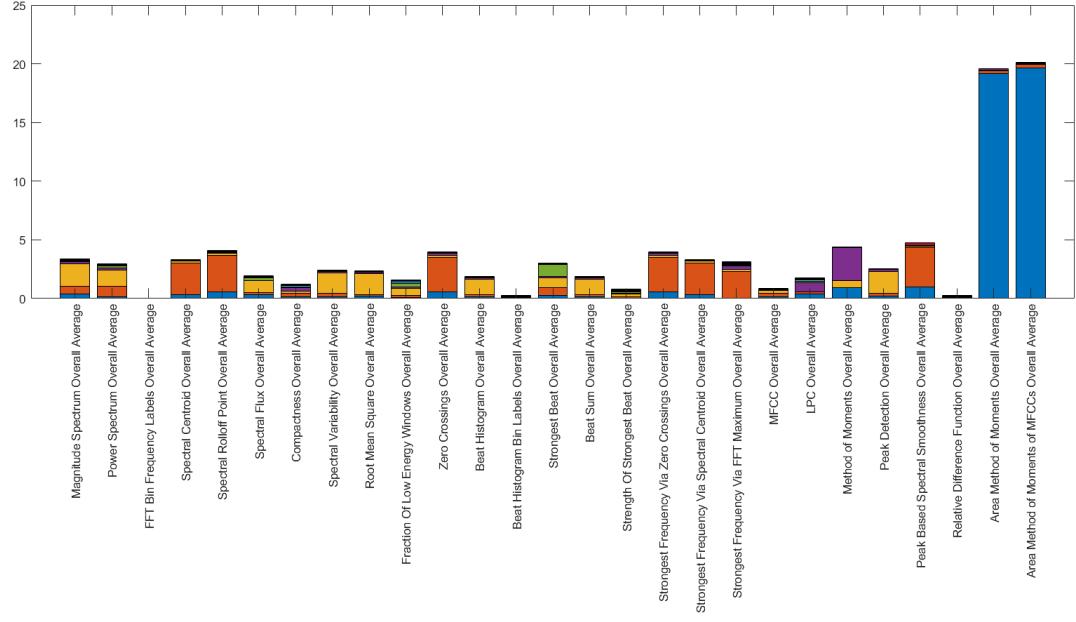


Figure 3.4: PCA coefficients weighted by eigen values (Normalized by Rescaling)

3.3 Scale Invariant Feature Transform (SIFT) Based Approach

Chapter 4

Implementation

Chapter 5

Results and Evaluation

5.1 Experiments

Chapter 6

Conclusions

References

- [1] E. D. N. W. Senevirathna and K. L. Jayaratne, “Radio Broadcast Monitoring to Ensure Copyright Ownership,” *International Journal on Advances in ICT for Emerging Regions (ICTer)*, vol. 11, p. 1, Aug. 2018.
- [2] Parliament of the democratic socialist republic of Sri Lanka, “Intellectual Property Act, No.36 of 2003.”
- [3] J. Serrà, E. Gómez, and P. Herrera, “Audio Cover Song Identification and Similarity: Background, Approaches, Evaluation, and Beyond,” in *Advances in Music Information Retrieval* (J. Kacprzyk, Z. W. Raś, and A. A. Wieczorkowska, eds.), vol. 274, pp. 307–332, Berlin, Heidelberg: Springer Berlin Heidelberg, 2010.
- [4] N. Kehtarnavaz, “Frequency Domain Processing,” *Digital Signal Processing System Design*, vol. 1, pp. 175–196, 2008.
- [5] Y. Ke, D. Hoiem, and R. Sukthankar, “Computer vision for music identification,” *Proceedings - 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, CVPR 2005*, vol. I, pp. 597–604, 2005.
- [6] D. G. Lowe, “Distinctive image features from scale-invariant keypoints,” *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [7] R. Sears, C. Van Ingen, and J. Gray, “To BLOB or Not To BLOB: Large Object Storage in a Database or a Filesystem?,” tech. rep., 2006.
- [8] M. Muja and D. G. Lowe, “Fast Approximate Nearest Neighbors with Automatic Algorithm Configuration,” tech. rep.