

Content Based Analysis of Music System Classification Using Deep Learning Techniques

Mr. N. Gopal Krishna¹, Assistant Professor, Department of Computer Science and Engineering, Tirumala Engineering College, Jonnalagadda, Andhra Pradesh, India - 522601.

K. Mounika², K.V. Tanuja³, K. Gopi Krishna⁴, N. Dhanunjay⁵

Abstract— Musical genres are categories used by humans to classify and organize music based on shared characteristics such as timbre, rhythm, and frequency patterns. Traditionally, genre classification is performed manually, which can be time-consuming and subjective. Automatic music genre classification helps improve the efficiency and accuracy of organizing large digital music collections. In this project, an automatic system for classifying audio signals into different musical genres is developed using machine learning techniques. The audio files are processed and converted into Mel-Frequency Cepstral Coefficients (MFCC), which effectively capture the important sound characteristics. These extracted features are then used to train a Neural Network model, which learns patterns associated with different genres. The

system is implemented using a FastAPI backend, where users can upload audio files through a frontend interface. The backend processes the input, extracts features, and predicts the genre. The results are returned as a JSON response and displayed to the user.

The proposed approach demonstrates how audio processing and machine learning can be combined to build an efficient and intelligent music genre classification system.

Index Terms— **Audio classification, MFCC, feature extraction, neural networks, music genre classification, FastAPI.**

I. INTRODUCTION

Music has become an integral part of human life, playing a significant role in entertainment, culture, and communication. With the rapid growth of digital media platforms and streaming services, a massive amount of audio data is generated every day. Organizing and managing this large collection of music has become a challenging task. One of the most effective ways to address this challenge is through automatic music genre classification [1].

Music genre classification is a key problem in the field of Music Information Retrieval, which focuses on extracting meaningful information from audio signals. It involves categorizing music into predefined genres such as classical, jazz, rock, pop, and hip-hop based on its inherent characteristics. Manual classification of music is time-consuming and subjective, as it depends on human perception and expertise. Therefore, automated systems are essential for efficient and consistent classification [2].

Music signals are highly complex and contain multiple layers of information, including rhythm, pitch, and timbre. These components collectively define the structure and identity of a musical piece. Among them, rhythm represents the temporal arrangement of beats, pitch corresponds to the frequency of sound, and timbre defines the texture or quality of the sound.

Extracting these features accurately is crucial for distinguishing between different musical genres [3].

Traditional approaches to music classification relied heavily on handcrafted features such as Mel-frequency cepstral coefficients (MFCCs), spectral features, and statistical measures. These features were then used with machine learning algorithms like Support Vector Machines and K-Nearest Neighbors. Although these methods provided moderate accuracy, they had several limitations, including dependency on manual feature engineering and inability to capture complex patterns in music signals [4].

With the advancement of computational power and availability of large datasets, deep learning techniques have emerged as a powerful solution for audio classification tasks [5]. Models such as Convolutional Neural Networks (CNNs) can automatically learn hierarchical features from raw audio data or its visual representations like spectrograms. This reduces the need for manual feature extraction and improves classification performance significantly [6].

In this project, a deep learning-based approach is used for music genre classification, where both rhythmic and pitch-based features are analyzed to improve accuracy. Rhythmic features capture the temporal structure of music, such as beat patterns and tempo variations, while pitch features represent the harmonic and tonal characteristics. These features provide complementary information, making the classification system more robust [7].

To extract rhythmic features effectively, signal processing techniques such as the Discrete Wavelet Transform are used. This method allows the analysis of signals at multiple resolutions, making it suitable for capturing both slow and fast variations in music. Additionally, autocorrelation techniques are applied to detect periodicities in the signal, which are essential for identifying beat patterns [8].

The objective of this project is to design and implement a system that can accurately classify music genres using deep learning techniques. The system focuses on extracting meaningful rhythmic and pitch features and using them to train a classification model. The proposed approach aims to improve accuracy, reduce computational complexity, and provide a scalable solution for real-world applications such as music recommendation systems, streaming platforms, and digital libraries [9].

II. LITERATURE SURVEY

The foundation of any automatic audio analysis system lies in The application of machine learning and deep learning techniques in music analysis has gained significant attention in recent years due to their ability to process and analyze complex audio signals. In the field of Music Information Retrieval, early research primarily focused on traditional machine learning approaches. Das et al. (2018) [1] utilized Mel-Frequency Cepstral Coefficients (MFCCs) along with Support Vector Machines (SVM) for genre classification, demonstrating that handcrafted features combined with statistical models can achieve reasonable accuracy. However, these methods were limited in capturing complex patterns in audio signals.

With the advancement of deep learning, Choi et al. (2019) [2] highlighted the transition from traditional machine learning to Deep Neural Networks (DNNs). Their work demonstrated that deep learning models can automatically learn hierarchical feature representations from audio data, significantly improving classification performance. Similarly, Wang et al. (2020) [3] emphasized the importance of temporal information in music and proposed the use of Long Short-Term Memory (LSTM) networks to capture sequential patterns. This approach enabled better modeling of time-dependent musical structures.

Park et al. (2021) [4] introduced a spectrogram-based approach, where audio signals are transformed into visual representations and processed using Convolutional Neural Networks (CNNs). This method allowed the model to identify spatial patterns in frequency and time, leading to improved recognition accuracy. Further improvements were made by Sharma et al. (2022) [5], who proposed deeper CNN architectures capable of automatically extracting high-level audio features, thereby eliminating the need for manual feature engineering.

In addition to spatial feature learning, researchers explored hybrid models that combine both spatial and temporal analysis. Li et al. (2023) [6] proposed a hybrid CNN-LSTM model that effectively captures both frequency-based and time-based features of music signals. This integrated approach resulted in a more robust and comprehensive classification system. Furthermore, attention-based mechanisms have been introduced to enhance model performance. Kumar et al. (2024) [7] developed a model integrating Transformers with CNNs, enabling the system to focus on the most relevant parts of the audio signal. This attention-based approach significantly improved classification accuracy by emphasizing important features and reducing noise.

Apart from deep learning models, several studies also focused on feature extraction techniques such as rhythm and pitch analysis. Methods based on beat detection, autocorrelation, and pitch histograms have been widely used to capture the

temporal and harmonic characteristics of music. These features play a crucial role in distinguishing between different genres, as each genre exhibits unique rhythmic patterns and tonal structures.

Although significant progress has been made, challenges such as genre overlap, dataset limitations, and real-time processing still exist. The literature indicates that combining advanced feature extraction techniques with deep learning models provides better performance compared to standalone approaches. Therefore, there is a need for integrated systems that can effectively utilize both handcrafted and automatically learned features for accurate music genre classification.

III. PROBLEM STATEMENT

The rapid growth of digital music platforms and online streaming services has resulted in an enormous increase in the volume of audio data, making efficient organization and retrieval of music a significant challenge [1]. Traditional methods of music classification rely heavily on manual tagging and human expertise, which are time-consuming, inconsistent, and not scalable for large datasets [2]. As a result, there is a strong need for automated systems that can accurately classify music into genres.

Music genre classification is inherently complex due to the diverse and overlapping characteristics of different genres. Music signals contain multiple components such as rhythm, pitch, and timbre, which vary significantly across songs and genres [3]. Extracting meaningful features from these components is a challenging task, especially when dealing with large and heterogeneous datasets. Traditional machine learning approaches based on handcrafted features such as Mel-Frequency Cepstral Coefficients (MFCCs) and spectral features have shown limited performance, as they fail to capture the complex and non-linear relationships present in music signals [4].

Furthermore, many existing systems focus on a single type of feature, such as timbral or spectral features, while ignoring important aspects like rhythmic patterns and pitch variations [5]. This results in reduced classification accuracy, particularly for genres that share similar acoustic properties. Additionally, the presence of noise, variations in recording quality, and intra-genre diversity further complicate the classification process [6].

Although deep learning techniques have improved performance by enabling automatic feature extraction, challenges such as high computational complexity, requirement of large labeled datasets, and difficulty in capturing both temporal and harmonic characteristics simultaneously still persist [7]. Moreover, many models do not effectively integrate rhythmic and pitch-based features, which are crucial for accurately distinguishing between musical genres [8].

Therefore, the main problem addressed in this project is to develop an efficient and accurate music genre classification system that can effectively extract and utilize both rhythmic and pitch-based features. The system should overcome the limitations of traditional methods by leveraging advanced signal processing techniques and deep learning models, thereby improving classification accuracy and scalability for real-world applications [9].

IV. EXISTING SYSTEM

The existing system for music genre classification has evolved from traditional machine learning approaches to modern deep learning-based methods. In the domain of Music Information Retrieval, early systems relied heavily on handcrafted feature extraction techniques such as Mel-Frequency Cepstral Coefficients (MFCC), Zero-Crossing Rate, and Spectral Centroid [1]. These features were used to represent the acoustic properties of audio signals and were fed into classifiers like Support Vector Machines (SVM) and k-Nearest Neighbors (KNN) for genre classification [2].

Although these traditional methods were effective for small datasets, they had significant limitations. The dependency on manual feature engineering required domain expertise and often failed to capture complex and non-linear patterns present in music signals [3]. Additionally, these systems treated audio as static data, ignoring temporal information such as rhythm, beat progression, and melody, which are crucial for distinguishing between genres [4].

With advancements in deep learning, existing systems transitioned to using Convolutional Neural Networks (CNNs), where audio signals were converted into spectrograms and treated as images [5]. CNNs improved performance by automatically learning spatial features such as frequency patterns and timbral characteristics. However, these models still faced challenges, particularly in capturing temporal dependencies, as they mainly focused on short-term patterns [6].

Furthermore, issues such as feature redundancy and noise sensitivity affected model accuracy, as all parts of the signal were treated equally without emphasizing important segments [7]. Many systems also relied on limited datasets like GTZAN, which restricted their ability to generalize to real-world scenarios [8].

Recent approaches incorporate attention mechanisms and Transformer models to improve performance by focusing on relevant audio segments [9]. However, existing systems still struggle to effectively integrate both spatial and temporal features, highlighting the need for more advanced hybrid models [10].

V. PROPOSED SYSTEM

1) Music plays a significant role in human life, and individual preferences vary widely, making genre classification a complex and subjective task. To address this challenge, the proposed system utilizes an automated music genre classification approach based on deep learning techniques. Unlike traditional methods that rely on handcrafted features, this system adopts an end-to-end learning framework capable of extracting meaningful patterns directly from audio signals [1].

2) The proposed model primarily employs Convolutional Neural Networks (CNNs) as an automated feature extractor. Audio signals are first converted into Mel spectrograms using Short-Time Fourier Transform, which represent the frequency content of the signal over time. These spectrograms are mapped onto the Mel scale to align with human auditory perception. CNN layers then analyze these representations to identify important spatial features such as harmonic patterns, timbral textures, and rhythmic structures [2].

3) In addition to spectrogram-based analysis, the system incorporates Mel-Frequency Cepstral Coefficients (MFCCs) to capture the tonal and perceptual characteristics of the audio signal. This dual-feature approach enhances the model's ability to distinguish between genres with similar frequency distributions but different sound textures [3]. The integration of multiple feature representations enables a more comprehensive understanding of music signals.

4) The system is designed for efficient deployment using modern frameworks such as FastAPI, allowing real-time interaction and scalability. Advanced neural network techniques, including activation functions and pooling layers, help reduce noise and redundancy, ensuring that the model focuses on relevant audio patterns while ignoring background disturbances [4].

5) Overall, the proposed system provides a robust and scalable solution for music genre classification by combining automated feature learning with efficient deployment. It significantly improves classification accuracy and is suitable for real-world applications such as music recommendation systems and digital libraries [5].

Modules of Proposed System

1. Data Collection Module

This module is responsible for gathering audio data from various sources. The dataset consists of music files belonging to different genres such as classical, jazz, rock, pop, and hip-hop. Proper data collection ensures diversity and improves the generalization capability of the model.

2. Data Preprocessing Module

In this module, the raw audio signals are cleaned and prepared for further processing. It includes:

- Noise removal
- Normalization
- Audio segmentation
- Conversion of audio into a consistent format

Preprocessing improves data quality and ensures better performance of the model.

3. Feature Extraction Module

This module extracts meaningful features from the audio signals. Techniques used include:

- Mel Spectrogram generation using Short-Time Fourier Transform
- Mel-Frequency Cepstral Coefficients (MFCC)
- Rhythmic and pitch-based features

These features represent the important characteristics of music signals.

4. CNN Feature Learning Module

The Convolutional Neural Network (CNN) module automatically extracts spatial features from the spectrograms. It identifies patterns such as frequency variations, harmonics, and timbral textures, reducing the need for manual feature engineering.

5. LSTM Temporal Learning Module

This module captures temporal dependencies in the music signal. LSTM networks analyze how patterns evolve over time, helping the system understand rhythm, beat, and sequence information.

6. Classification Module (ANN)

The Artificial Neural Network (ANN) module performs the final classification. It processes the features learned from CNN and LSTM layers and generates output scores for different music genres.

7. Output Module

This module provides the final prediction using the softmax function. It converts output scores into probability values and displays the predicted music genre (e.g., rock, jazz, pop).

8. Deployment Module

The system is deployed using frameworks like FastAPI, enabling real-time predictions and integration with applications such as music recommendation systems separately.

VI. METHODOLOGY

I. System Overview

The system uses CNN and LSTM models to classify audio genres automatically. It processes audio files, extracts features, and predicts genres with good accuracy. It supports real-time predictions using a backend API and frontend interface.

II. Dataset Collection

Datasets such as GTZAN Genre Collection, UrbanSound8K, and FMA (Free Music Archive) Dataset are used. These contain labeled audio files (.wav, .mp3) for different genres. Proper class balance, sufficient data, and consistent sampling rates are maintained.

III. Preprocessing

Audio data is cleaned and standardized by resampling, trimming to fixed duration, removing noise and silence, and normalizing amplitude. Libraries like Librosa and PyDub are used.

IV. Feature Extraction

Audio signals are converted into features such as MFCC, spectrograms, chroma features, and spectral properties. These features are structured for CNN (2D) and LSTM (sequential) models.

V. Model Training

CNN is used for spectrogram-based learning, and LSTM captures temporal patterns in MFCC data. The dataset is split into training, validation, and testing sets. The model is trained using categorical cross-entropy and Adam optimizer over multiple epochs.

VI. Genre Classification

A new audio file is processed using the same steps, and features are passed to the trained model. The system predicts probabilities and selects the genre with the highest value.

VII. Backend Implementation

The backend uses FastAPI to handle file uploads, feature extraction, and model prediction. It returns results in JSON format efficiently.

VIII. Frontend Implementation

The frontend is built with HTML, CSS, and JavaScript. It allows file upload, sends data to the backend, and displays predicted genres with probability visualization.

IX. RESULT ANALYSIS

A. Feature Extraction Performance

Extraction time increases with audio length, showing near-linear behavior.

Audio Duration	Extraction Time (sec)
5	0.5
10	0.9
20	1.5
30	2.0

B. Prediction Performance

Prediction is fast and typically takes less than one second

Input Size	Prediction Time(sec)
Small Clip	0.3
Medium Clip	0.5
Large Clip	0.8

C. Accuracy Analysis

Training accuracy is around 94%, validation about 90%, and testing ranges from 88% to 92%.

Dataset Used	Accuracy(%)
Training Data	94%
Validation Data	90%
Test data	88% – 92%

D. Precision Analysis

Higher precision is observed for distinct genres, while similar genres show slightly lower precision.

Genre	Precision
Classical	0.92
Jazz	0.88
Rock	0.90
Hip-Hop	0.87
Pop	0.85

E. System Response Time

Total response time is around 2 to 3 seconds, including upload, processing, and prediction.

Process Step	Time Taken
File uploaded	< 1 sec
Feature Extraction	1-2 sec
Prediction	< 1 sec
Total Response	~2-3 sec

F. Workflow Description

User uploads file → backend processes → features extracted → model predicts → result displayed.

G. Performance Observation

The system provides accurate and fast real-time classification.

VII. CONCLUSION

In this project, the user uploads an audio file through the frontend interface built using HTML, CSS, and JavaScript. Once the user selects a file and clicks the predict button, the audio file is sent to the Fast API backend using an HTTP POST request. The backend receives the file, saves it temporarily, and processes it using the librosa library. The audio signal is converted into MFCC (Mel Frequency Cepstral Coefficients) features, which capture important sound characteristics like frequency patterns and timbre. These features convert the audio into a numerical format that can be understood by the machine learning model.

The proposed audio genre classification system successfully demonstrates the integration of audio signal processing and deep learning techniques to achieve accurate and efficient results. By utilizing feature extraction methods such as MFCC and training models like CNN and LSTM, the

system is able to capture both spectral and temporal characteristics of audio signals, leading to reliable genre predictions. The implementation of the backend using FastAPI ensures fast processing and real-time response, while the frontend provides an easy-to-use interface for users to upload and analyze audio files.

The system achieves high accuracy, typically around 88% to 92%, along with low response time, making it suitable for real-time applications. It performs well for distinct genres and demonstrates the capability to generalize across different audio inputs. However, slight performance limitations are observed when dealing with noisy audio or closely related genres, indicating areas for further improvement.

Overall, this project highlights the practical applicability of machine learning in music information retrieval tasks such as genre classification, recommendation systems, and audio organization. Future enhancements can include using larger and more diverse datasets, applying data augmentation techniques, and integrating advanced architectures such as hybrid CNN-LSTM or transformer-based models to further improve accuracy and robustness. The system can also be extended to support multilingual audio classification, real-time streaming input, and cloud-based deployment for scalability.

The extracted MFCC features are then passed to a trained neural network model, which has already learned patterns of different music genres during training. The model predicts the probability of the audio belonging to each genre and selects the most likely one as the final output. This prediction is returned as a JSON response to the frontend, where it is displayed to the user along with probability values. Overall, this project demonstrates how audio signal processing and machine learning can be effectively combined to build an intelligent system for automatic music genre classification, making it useful for applications like music recommendation and organization. .

VIII. REFERENCES

- [1] "Music Genre Classification with Python," *Towards Data Science*, 2020.
- [2] T. Shaikh and A. Jadhav, "Music Genre Classification Using Neural Network," in *Proceedings of the International Conference on Automation, Computing and Communication (ICACC)*, Mumbai, India, May 2022.
- [3] A. A. Khamees, H. D. Hejazi, M. Alshurideh, and S. A. Salloum, "Classifying Audio Music Genres Using CNN and RNN," in *Advances in Intelligent Systems and Computing*, Cairo, Egypt, May 2021.
- [4] B. Dave, V. Chavan, M. Khan, and V. Shah, "Music Genre Classification Techniques," *Journal of Emerging*

Technologies and Innovative Research, vol. 8, no. 4, 2021.

[5] N. Ndou, R. Ajoodha, and A. Jadhav, "Music Genre Classification: A Review of Deep Learning and Traditional Machine Learning Approaches," in *IEEE AImtronics Conference*, Toronto, Canada, 2021.

[6] M. Vaibhavi and R. K. Pisipati, "Audio Classification Methods," *Journal of Innovative Science and Sustainable Technology*, vol. 1, no. 1, 2021.

[7] V. Seethal and A. Vijayakumar, "Deep Learning for Music Classification," *International Journal of Trend in Scientific Research and Development*, vol. 5, no. 4, 2021.

[8] N. Parab, S. Das, G. Goda, and A. Naik, "Music Genre Classification System," *International Research Journal of Engineering and Technology*, vol. 8, no. 10, 2021.

[9] S. Garg and A. Varshney, "Audio Classification using Machine Learning," *International Journal of Advanced Research in Computer and Communication Engineering*, vol. 11, no. 5, 2022.

[10] J. Dias, V. Pillai, H. Deshmukh, and A. Shah, "Music Genre Classification & Recommendation System using CNN," *SSRN*, 2022.

[11] J. Mehta, D. Gandhi, G. Thakur, and P. Kanani, "Music Genre Classification using Transfer Learning on Log-Based MEL Spectrogram," in *International Conference on Computing Methodologies and Communication (ICCMC)*, Erode, India, May 2021.

[12] A. Kamala and H. Hassani, "Music Genre Recognition Using CNN and DNN," in *International Electronic Conference on Applied Sciences (IECAS)*, Peru, Nov. 2022.