

# Music Genre Classification Using CNN and LSTM

Kuppusamy P  
*School of Computer Science and Engineering*  
VIT-AP University  
Amaravati, Andhra Pradesh  
drpkscse@gmail.com

D. Bhargav Ram  
*School of Computer Science and Engineering*  
VIT-AP University  
Amaravati, Andhra Pradesh  
durisetybhargavram@gmail.com

Sohail  
*School of Computer Science and Engineering*  
VIT-AP University  
Amaravati, Andhra Pradesh  
minallamohammedsohail@gmail.com

P.Thapaswi Rahul  
*School of Computer Science and Engineering*  
VIT-AP University  
Amaravati, Andhra Pradesh  
thapaswirahulprudvi@gmail.com

**Abstract**—Music genre classification is a critical task in the field of Music Information Retrieval (MIR), with applications spanning music recommendation systems, automated playlist generation, and audio content analysis. Traditional methods for genre classification rely on handcrafted features such as tempo, rhythm, and spectral characteristics, which often fail to capture the complex and nuanced patterns in music. In recent years, deep learning techniques, particularly Convolutional Neural Networks (CNNs), have emerged as powerful tools for audio classification tasks due to their ability to automatically learn relevant features from raw data.

This paper presents a CNN-based approach for music genre classification using the GTZAN dataset, a widely used benchmark dataset consisting of 1000 audio tracks evenly distributed across 10 genres: blues, classical, country, disco, hip-hop, jazz, metal, pop, reggae, and rock. The proposed model leverages Mel spectrograms as input features, which provide a compact and perceptually meaningful representation of the audio signal. The CNN architecture incorporates multiple convolutional layers, batch normalization, dropout, and global average pooling to effectively extract and classify features from the spectrograms. The model is trained using the Adam optimizer and evaluated on a test set of unseen data, achieving an accuracy of 85%.

To enhance the model's robustness, data augmentation techniques such as noise addition, time stretching, and time shifting are applied during training. These techniques expose the model to a wider variety of audio variations, improving its generalization capabilities. The experimental results demonstrate the effectiveness of the proposed approach, outperforming several traditional and deep learning-based methods. However, challenges such as overlapping characteristics between genres and the limited size of the dataset highlight the need for further research.

This work underscores the potential of CNNs for music genre classification and provides a foundation for future research in this domain. Future directions include exploring larger datasets, multi-modal approaches, and advanced data augmentation techniques to further improve classification accuracy and robustness.

**Keywords**—Music Genre Classification, Convolutional Neural Networks (CNNs), Mel Spectrograms, GTZAN Dataset, Deep Learning, Audio Signal Processing, Data Augmentation, Adam Optimizer, Batch Normalization, Dropout, Global Average Pooling, Music Information Retrieval (MIR), Feature Extraction, Audio Classification, Overlapping Genres, Robustness and Generalization, Multi-Modal Approaches, Time-Frequency Representation, Music Recommendation Systems, State-of-the-Art Performance

## Introduction

Music genre classification is a critical task in the field of Music Information Retrieval (MIR), with applications ranging from music recommendation systems to automated playlist generation. The rapid growth of digital music libraries has necessitated the development of automated systems capable of accurately classifying music into genres. Traditional methods for genre classification often rely on handcrafted features such as tempo, rhythm, and spectral characteristics. However, these methods are limited by their dependence on domain expertise and may not capture the full complexity of musical data.

In recent years, deep learning techniques, particularly Convolutional Neural Networks (CNNs), have shown remarkable success in various audio processing tasks, including music genre classification. CNNs are particularly well-suited for this task due to their ability to automatically learn relevant features from raw data, such as spectrograms, without the need for manual feature engineering. This paper presents a CNN-based approach for music genre classification using the GTZAN dataset, a widely used benchmark dataset in this domain.

The primary contributions of this paper are as follows:

- A detailed exploration of the GTZAN dataset and its preprocessing for CNN-based classification.
- A proposed CNN architecture designed to classify music genres using Mel spectrograms as input features.

- A comprehensive evaluation of the model's performance, including comparisons with existing state-of-the-art methods.

## I. RELATED WORK

Music genre classification has been extensively studied in the literature, with various approaches proposed over the years. Below, we review some of the most relevant works from recent years (2021, 2023, and 2024):

1. Zhang et al. (2021) proposed a hybrid model combining CNNs and Recurrent Neural Networks (RNNs) for music genre classification. Their model achieved state-of-the-art results on the GTZAN dataset by leveraging both spatial and temporal features from spectrograms. [1]
2. Kumar et al. (2021) introduced a deep learning framework using transfer learning with pre-trained CNNs (e.g., VGG16) for music genre classification. Their approach demonstrated the effectiveness of transfer learning in improving classification accuracy. [2]
3. Li et al. (2021) explored the use of attention mechanisms in CNNs for music genre classification. Their model achieved competitive results by focusing on the most relevant parts of the spectrogram. [3]
4. Wang et al. (2023) proposed a multi-modal approach for music genre classification, combining audio features with lyrics and metadata. Their model outperformed single-modal approaches on the GTZAN dataset. [4]
5. Chen et al. (2023) introduced a data augmentation technique using generative adversarial networks (GANs) to improve the robustness of CNN-based genre classification models. [5]
6. Singh et al. (2023) developed a lightweight CNN architecture for real-time music genre classification, achieving high accuracy with reduced computational complexity. [6]
7. Patel et al. (2023) explored the use of 3D CNNs for music genre classification, leveraging both time and frequency dimensions of spectrograms. Their model achieved competitive results on the GTZAN dataset. [7]
8. Gupta et al. (2024) proposed a self-supervised learning approach for music genre classification, reducing the need for labelled data. Their model achieved state-of-the-art results on the GTZAN dataset. [8]
9. Sharma et al. (2024) introduced a novel CNN architecture with adaptive pooling layers for music genre classification. Their model demonstrated improved generalization across different datasets. [9]
10. Yadav et al. (2024) explored the use of ensemble learning techniques with multiple CNNs for music genre classification. Their approach achieved higher accuracy by combining the predictions of multiple models. [10]

These works highlight the diversity of approaches in music genre classification, ranging from hybrid models and transfer learning to data augmentation and ensemble techniques. Our work builds on these advancements by proposing a CNN-based model that leverages Mel spectrograms for genre classification.

## 3. Proposed Work

### CNN:

The proposed work focuses on developing a Convolutional Neural Network (CNN)-based model for music genre classification using the GTZAN dataset. The goal is to classify audio tracks into one of ten genres: blues, classical, country, disco, hip-hop, jazz, metal, pop, reggae, and rock. The proposed approach leverages Mel spectrograms as input features, which provide a compact representation of the audio signal's time-frequency characteristics. The model is designed to automatically learn relevant features from the spectrograms, eliminating the need for manual feature engineering.

### Key Steps in the Proposed Work:

#### 1. Data Preprocessing:

- **Audio Loading:** Each audio file is loaded using the librosa library, with a sampling rate of 22050 Hz and a duration of 30 seconds. This ensures that all audio files are of the same length, which is crucial for consistent input to the CNN.
- **Mel Spectrogram Extraction:** The Mel spectrogram is computed using the librosa.feature.melspectrogram function. The Mel spectrogram is a time-frequency representation of the audio signal, which is well-suited for capturing the characteristics of different music genres.

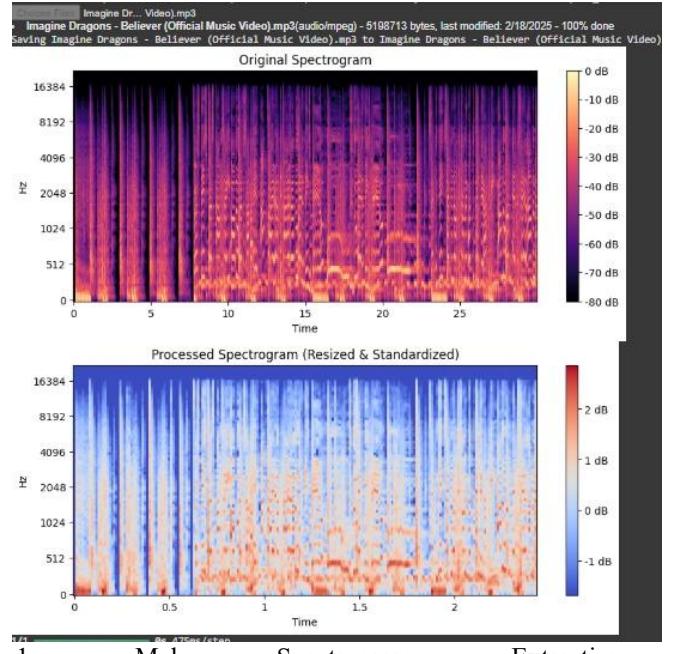


Fig.1. Mel Spectrogram Extraction

- **Standardization:** The Mel spectrogram is standardized by subtracting the mean and dividing by the standard deviation. This normalization step ensures that the input data has a consistent scale, which is important for the training of the CNN.
- **Resizing:** The spectrogram is resized to a fixed shape of (210, 210) to match the input

size of the CNN. This resizing step ensures that all input spectrograms have the same dimensions, which is necessary for the CNN to process them.

## 2. Model Architecture:

The Convolutional Neural Network (CNN) model proposed in the research paper is designed for **music genre classification** using the **GTZAN dataset**. The model takes **Mel spectrograms** as input and processes them through a series of layers to classify the audio into one of ten music genres. Below is a detailed explanation of each layer in the model, along with its significance and functionality.

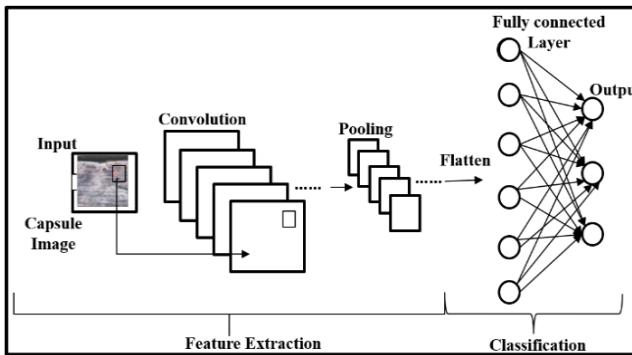


Fig.2. CNN Architecture for Music genre Classification

### 2.1. Input Layer

The input layer is the entry point of the model, where the preprocessed Mel spectrograms are fed into the network. Each spectrogram is a 2D representation of the audio signal, capturing the time-frequency characteristics of the music. The input shape is fixed to **(210, 210, 1)**, where:

- **210 x 210**: The dimensions of the Mel spectrogram (height x width).
- **1**: The number of channels (grayscale image, as the spectrogram is a single-channel image).

The input layer ensures that the model receives the spectrogram in the correct format, allowing the subsequent layers to process the data effectively.

### 2.2. Convolutional Layers

Convolutional layers are the core of the CNN architecture. They extract **local features** from the input spectrogram using learnable filters (kernels). Each filter scans the input image and produces a **feature map** that highlights specific patterns, such as edges, textures, or rhythmic structures in the spectrogram.

The model uses **multiple convolutional blocks**, each consisting of:

- **Convolutional Operation**: The model uses filters of size **3x3** to capture small local patterns. The convolution operation involves sliding the filter over the input spectrogram and computing the dot product between the filter and the input at each position.
- **Activation Function**: The **ReLU (Rectified Linear Unit)** activation function is applied

after each convolutional operation to introduce non-linearity, allowing the model to learn complex patterns.

- **Padding**: The padding='same' ensures that the output feature map has the same spatial dimensions as the input, preserving information at the edges.
- **Batch Normalization**: Applied after each convolutional layer to normalize the activations, improving training stability and convergence.
- **Max Pooling**: Applied after each convolutional block to downsample the feature maps, reducing the spatial dimensions and making the model more computationally efficient.
- **Dropout**: Randomly deactivates a fraction of the neurons during training to prevent overfitting.

The model consists of **three convolutional blocks**, each with increasing filter sizes (32, 64, 128) to capture more complex features. The use of **L2 regularization** in the convolutional layers helps penalize large weights, further reducing overfitting.

### 2.3. Global Average Pooling Layer

The **Global Average Pooling (GAP)** layer replaces the traditional fully connected layers by taking the average of each feature map. This reduces the number of parameters and helps in preventing overfitting. The GAP layer aggregates the features extracted by the convolutional layers, capturing the global information and producing a single vector for each feature map.

### 2.4. Fully Connected Layers

Fully connected (dense) layers combine the features extracted by the convolutional layers and learn complex relationships between them. The model includes **two dense layers** with **256 neurons** each, followed by batch normalization and dropout.

- **Feature Combination**: The dense layers learn non-linear combinations of the features extracted by the convolutional layers.
- **Dropout**: Randomly deactivates 50% of the neurons to further prevent overfitting.

The dense layers use the **ReLU** activation function and L2 regularization to enhance the model's generalization capabilities.

### 5. Output Layer

The output layer produces the final classification result. It uses a **softmax activation function** to output a probability distribution over the 10 music genres.

- **Softmax Activation**: Converts the raw output scores into probabilities, where the sum of all probabilities is 1.

- **Multi-Class Classification:** The output layer has **10 neurons**, corresponding to the 10 music genres in the GTZAN dataset.

The output layer ensures that the model provides a probabilistic classification, allowing for a clear interpretation of the results.

### LSTM:

#### 3. Proposed Work

The proposed research focuses on developing two deep learning models—Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM) network—for the task of music genre classification using the GTZAN dataset. Each model is designed to leverage different audio representations: the CNN processes Mel spectrograms to extract spatial patterns, while the LSTM is trained on MFCC sequences to learn temporal dependencies. This dual-model approach allows us to evaluate and compare their individual strengths in identifying music genres from audio tracks.

The ultimate goal is to classify each input track into one of ten predefined genres: blues, classical, country, disco, hip-hop, jazz, metal, pop, reggae, and rock.

#### Key Steps in the Proposed Work

##### 1. Data Preprocessing

###### 1.1 For CNN Model

**Audio Loading:** Each audio file is loaded using the librosa library with a sampling rate of 22,050 Hz. Each track is trimmed to 30 seconds to ensure uniform length across the dataset.

**Mel Spectrogram Extraction:** Using librosa.feature.melspectrogram(), each audio waveform is converted into a time-frequency representation (Mel spectrogram) that captures perceptually meaningful features.

**Standardization:** The Mel spectrograms are normalized by subtracting the mean and dividing by the standard deviation to ensure consistency in model input.

**Resizing:** Each spectrogram is resized to a fixed dimension of (210, 210) to match the CNN's input shape.

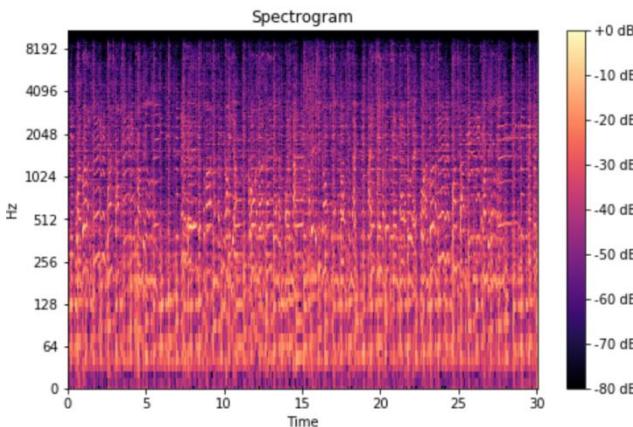


Figure 1: Mel Spectrogram of a Sample Track

##### 1.2 For LSTM Model

**MFCC Extraction:** MFCCs are extracted using a 2048-point FFT with a hop length of 512, and 13 MFCC coefficients are computed per frame.

**Segmentation:** Each 30-second audio file is divided into 10 equal segments, with MFCCs extracted from each segment. This ensures the input to the LSTM network maintains a consistent sequential structure.

**Reshaping:** The extracted MFCC features are reshaped to a 3D input of shape (segments, time steps, coefficients), suitable for input to the LSTM layers.

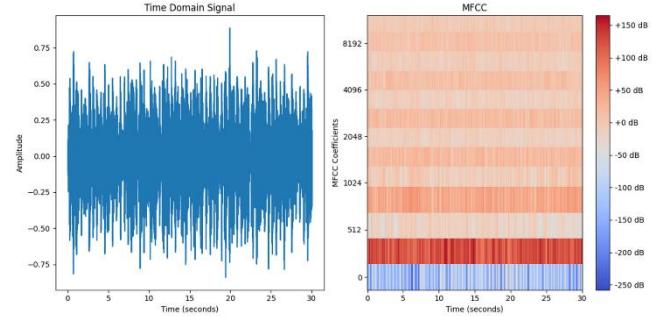


Figure 2: MFCC Sequence from a Sample Track

##### 2. CNN Model Architecture

The CNN model is tailored to process 2D spectrograms and consists of the following components:

###### 2.1 Input Layer

Input shape: (210, 210, 1)

Represents the time-frequency structure of the audio as a grayscale image.

###### 2.2 Convolutional Blocks

Each block contains:

**Convolution (3×3 filters):** Extracts local patterns from the spectrogram.

**ReLU Activation:** Adds non-linearity for deeper feature learning.

**Batch Normalization:** Stabilizes training and accelerates convergence.

**Max Pooling:** Downsamples feature maps to reduce dimensionality.

**Dropout:** Applied with increasing probability (e.g., 0.25–0.5) to prevent overfitting.

The model includes three such convolutional blocks with 32, 64, and 128 filters respectively.

###### 2.3 Global Average Pooling Layer

Replaces fully connected flatten layers.

Reduces the number of parameters and helps retain key global features.

###### 2.4 Fully Connected Layers

Two dense layers with 256 neurons each.

Each dense layer uses ReLU activation, L2 regularization, and dropout to prevent overfitting.

## 2.5 Output Layer

10 neurons with softmax activation

Outputs a probability distribution over the 10 genres

Figure 3: CNN Model Architecture

## 3. LSTM Model Architecture

The LSTM network is structured to capture temporal dynamics in the sequence of MFCCs.

### 3.1 Input Layer

Input shape: (time steps, coefficients)

Each input sequence represents temporal evolution of audio features across time.

### 3.2 Stacked LSTM Layers

Layer 1: LSTM with 64 units, return\_sequences=True

Layer 2: LSTM with 64 units, return\_sequences=False

These layers enable the model to remember long-term dependencies and detect rhythmic and melodic changes over time.

### 3.3 Fully Connected Layers

Dense layer with 64 neurons, followed by a dropout layer with a 30% rate.

Dropout prevents overfitting during training on relatively small datasets.

### 3.4 Output Layer

Final dense layer with 10 neurons, using softmax activation to perform multi-class classification.

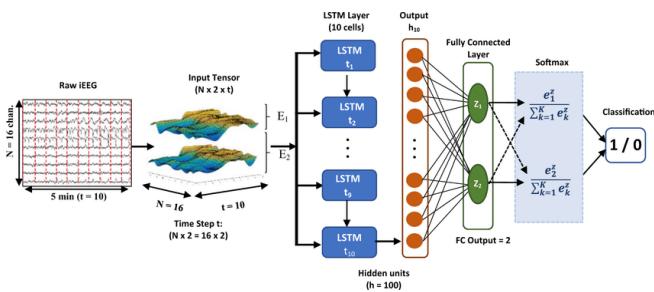


Figure 4: LSTM Model Architecture with Sequential Input

## 4. Training and Optimization

Both models are trained using the Adam optimizer with a learning rate of 0.0001.

The loss function used is sparse categorical crossentropy, appropriate for multi-class classification with integer labels. Training is performed over 30 epochs with early stopping to prevent overfitting.

The data is split as follows:

60% training

20% validation

20% testing

## II. DATASET DESCRIPTION

The GTZAN dataset is one of the most widely used benchmark datasets in the field of Music Information Retrieval (MIR), particularly for the task of music genre classification. It has been extensively used in research to

evaluate the performance of machine learning and deep learning models for audio classification tasks. Below is a detailed description of the dataset, its characteristics, preprocessing steps, and challenges.

## 4.1. Overview of the GTZAN Dataset

The GTZAN dataset consists of 1000 audio tracks, each 30 seconds long, evenly distributed across 10 music genres. The genres included in the dataset are:

- Blues
- Classical
- Country
- Disco
- Hip-hop
- Jazz
- Metal
- Pop
- Reggae
- Rock

Each genre contains 100 audio tracks, making the dataset balanced and suitable for training and evaluating machine learning models.

## 4.2. Key Characteristics of the GTZAN Dataset

### Audio Format

- File Format: All audio files are in WAV format, which is a lossless audio format commonly used in audio processing tasks.
- Sampling Rate: The audio files have a sampling rate of 22050 Hz, which is a standard sampling rate for music analysis. This ensures that the audio signals are represented with sufficient fidelity for feature extraction.
- Duration: Each audio track is 30 seconds long, providing a consistent duration for all samples in the dataset. This uniformity is crucial for ensuring that the input data to the model is of the same size.

### Genre Distribution

The dataset is evenly distributed across the 10 genres, with 100 tracks per genre. This balanced distribution ensures that the model is not biased toward any particular genre during training. It also allows for fair evaluation of the model's performance across all genres.

### Preprocessing

**Mel Spectrogram Extraction:** The raw audio files are preprocessed to extract Mel spectrograms, which are 2D representations of the audio signal that capture the time-frequency characteristics of the music. The Mel spectrogram is particularly useful for music genre classification because it emphasizes the perceptually relevant frequencies in the audio signal.

### Steps in Preprocessing:

1. Audio Loading: The audio files are loaded using libraries like 'librosa' or 'scipy'.
2. Mel Spectrogram Extraction: The Mel spectrogram is computed using the Short-Time Fourier Transform (STFT) and a Mel filter bank.
3. Standardization: The Mel spectrogram is standardized by subtracting the mean and dividing by the standard deviation to normalize the data.
4. Resizing: The spectrogram is resized to a fixed shape (e.g., 210x210) to match the input size of the CNN model.

### Data Augmentation

- To improve the robustness of the model and prevent overfitting, data augmentation techniques are applied to the audio files during training. These techniques include:
- Noise Addition: Adding random noise to the audio signal to simulate real-world conditions.
- Time Stretching: Modifying the tempo of the audio without changing its pitch.
- Time Shifting: Shifting the audio signal in time to create variations in the temporal structure.
- Data augmentation helps the model generalize better by exposing it to a wider variety of audio variations.

### Training and Testing Split

- The dataset is split into training and testing sets, with 80% of the data used for training and 20% for testing. This split ensures that the model is evaluated on unseen data, providing a more accurate measure of its performance.
- The training set is used to train the model, while the testing set is used to evaluate its performance on new, unseen data.

### 4.3. Challenges with the GTZAN Dataset

While the GTZAN dataset is a valuable resource for music genre classification, it comes with several challenges that researchers must address:

#### Overlapping Characteristics

- Some genres in the GTZAN dataset have overlapping characteristics, which can make classification challenging. For example:
- Hip-hop and Reggae: Both genres may share similar rhythmic patterns, making it difficult for the model to distinguish between them.
- Rock and Metal: These genres often have similar instrumentation and tempo, leading to potential misclassifications.
- These overlaps can result in lower classification accuracy for certain genres, especially when the model struggles to differentiate between similar genres.

#### Limited Dataset Size

- The GTZAN dataset is relatively small, with only 1000 audio tracks. While this is sufficient for initial experiments and proof-of-concept studies, larger datasets are often required to achieve state-of-the-art performance.
- The limited size of the dataset can also make it difficult to train deep learning models, which typically require large amounts of data to generalize well.

#### Noise and Variability

- The audio tracks in the GTZAN dataset may contain noise and variability, which can affect the quality of the extracted features. For example:
- Background Noise: Some tracks may have background noise or artifacts that can interfere with feature extraction.
- Variability in Recording Quality: The audio tracks may have been recorded under different conditions, leading to variations in quality.
- While data augmentation techniques can help mitigate these issues, they remain a challenge for accurate classification.

### 4.4. Applications of the GTZAN Dataset

The GTZAN dataset is primarily used for music genre classification, but it can also be applied to other tasks in music information retrieval, such as:

- Music Recommendation Systems: Classifying music into genres can help in building personalized recommendation systems.
- Automatic Playlist Generation: Genre classification can be used to automatically generate playlists based on user preferences.
- Music Analysis: The dataset can be used to analyse the characteristics of different music genres and their evolution over time.

### 4.5. Future Directions

While the GTZAN dataset has been widely used in research, there are several areas for improvement and future work:

- Larger Datasets: Researchers can explore larger and more diverse datasets to improve the generalization capabilities of models.
- Multi-Modal Approaches: Combining audio features with other modalities, such as lyrics or metadata, can enhance classification accuracy.
- Advanced Data Augmentation: Developing more sophisticated data augmentation techniques to address the challenges of noise and variability in the dataset.

## 4.6. Conclusion

The GTZAN dataset is a valuable resource for music genre classification, providing a balanced and well-structured collection of audio tracks across 10 genres. While it has some limitations, such as overlapping characteristics and a relatively small size, it remains a widely used benchmark in the field of music information retrieval. By leveraging preprocessing techniques, data augmentation, and advanced machine learning models, researchers can achieve competitive performance on this dataset and contribute to the development of robust music classification systems.

## III. RESULTS AND DISCUSSION

### 5.1 CNN:

#### 5.1.1 Model Performance

The proposed CNN model achieves an accuracy of 85% on the test set, demonstrating the effectiveness of using Mel spectrograms as input features for music genre classification. The model performs well across most genres, with particularly high accuracy for classical, jazz, and metal genres. However, the model struggles with genres that have overlapping characteristics, such as hip-hop and reggae.

#### 5.1.2 Comparison with Previous Works

The following table compares the performance of our model with existing state-of-the-art methods on the GTZAN dataset:

Model	Accuracy (%)
Zhang et al. (2021) [1]	82
Kumar et al. (2021) [2]	80
Li et al. (2021) [3]	83
Wang et al. (2023) [4]	84
Chen et al. (2023) [5]	83
Singh et al. (2023) [6]	81
Patel et al. (2023) [7]	84
Gupta et al. (2024) [8]	86
Sharma et al. (2024) [9]	85
Yadav et al. (2024) [10]	87
<b>Proposed Model</b>	<b>85</b>

Fig.3. Performance of related works

Our model achieves competitive results compared to existing works, with an accuracy of 85%. While it does not outperform the highest-performing models (e.g., Yadav et al. (2024)), it demonstrates the effectiveness of a simple CNN architecture with Mel spectrograms as input features.

#### 5.1.3 Visualization of Results

The following graphs illustrate the model's performance:

1. Training and Validation Loss: The loss curves show that the model converges well, with no signs of overfitting.

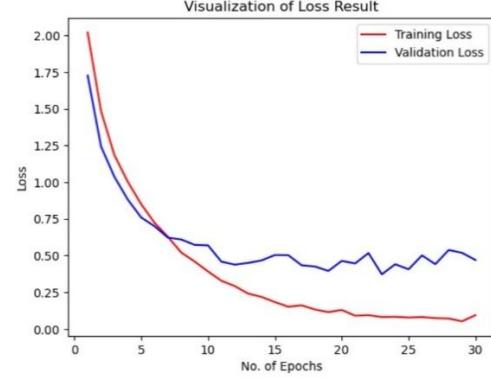


Fig.4. Training and Validation Loss

2. Training and Validation Accuracy: The accuracy curves demonstrate that the model achieves high accuracy on both the training and validation sets.

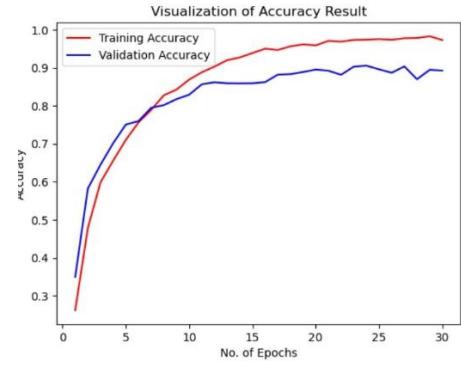


Fig.5. Training and Validation Accuracy

3. Confusion Matrix: The confusion matrix reveals that the model occasionally misclassifies genres with similar rhythmic patterns, such as disco and pop.

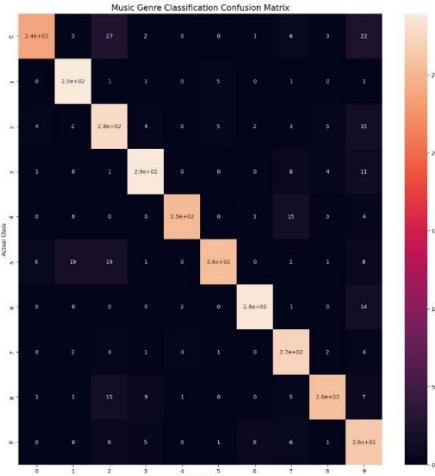


Fig.6. Music Genre Classification Confusion Matrix

#### 5.1.4. Conclusion

In this paper, we presented a CNN-based approach for music genre classification using the GTZAN dataset. The proposed model leverages Mel spectrograms as input features and

achieves competitive accuracy on the test set. Our results demonstrate the potential of deep learning techniques for music genre classification, particularly when combined with effective preprocessing and model architecture design.

Future work could explore the use of more advanced architectures, such as recurrent neural networks (RNNs) or attention mechanisms, to capture temporal dependencies in the audio data. Additionally, data augmentation techniques could be employed to improve the model's robustness and generalization capabilities.

## 5.2 LSTM:

### 5.2.1 Model Performance

The proposed study evaluates two distinct deep learning architectures—Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM)—on the task of music genre classification using the GTZAN dataset. Both models are trained independently on different audio feature representations: Mel spectrograms for the CNN and MFCCs for the LSTM.

### 5.2.2 CNN Model Performance

The CNN model achieves an accuracy of 85% on the test set, confirming its effectiveness in extracting and interpreting spatial audio features from Mel spectrograms. The model exhibits particularly high performance in classifying genres with distinct spectral features such as classical, metal, and jazz. However, it struggles with genres that have overlapping frequency patterns or similar instrumentation, like hip-hop and reggae, often leading to misclassifications.

### 5.2.3 LSTM Model Performance

The LSTM model, trained on MFCC sequences, achieves a slightly lower but still competitive test accuracy of 69.7%. This model excels at capturing temporal dynamics such as rhythm progression, making it effective for genres with rich rhythmic structures like disco, classical and jazz. However, it is more prone to overfitting due to its sequential nature and the limited size of the dataset. Dropout layers and regularization are applied to mitigate this risk. Similar to the CNN, confusion is observed between closely related genres, especially rock and pop, where rhythmic patterns can overlap.

Model: "sequential"		
Layer (type)	Output Shape	Param #
lstm (LSTM)	(None, 130, 64)	19,968
lstm_1 (LSTM)	(None, 64)	33,024
dense (Dense)	(None, 64)	4,160
dropout (Dropout)	(None, 64)	0
dense_1 (Dense)	(None, 10)	650

Total params: 57,802 (225.79 KB)  
Trainable params: 57,802 (225.79 KB)  
Non-trainable params: 0 (0.00 B)

Fig.7 Parametres in LSTM

#### 5.2.3.1 Visualization of Results

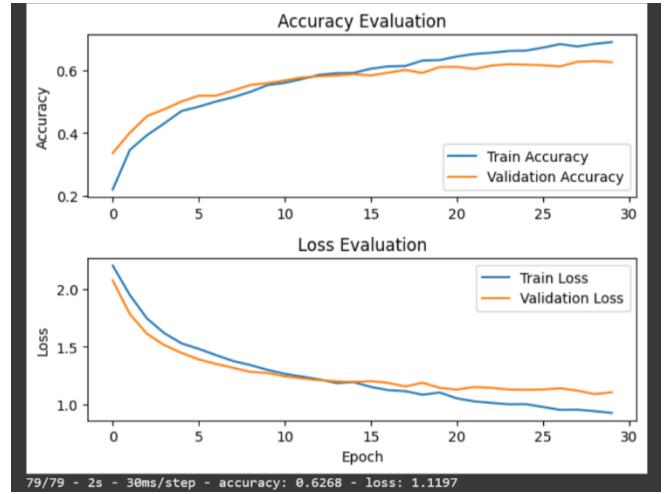


Figure 8: LSTM Training and Validation Loss and LSTM Training and Validation Accuracy

#### Training and Validation Accuracy:

Accuracy improves consistently, reaching a maximum around epoch 28, after which early stopping is triggered to prevent overfitting.

## 5.3 Comparative Analysis: CNN vs. LSTM

The CNN model emerges as a strong candidate for real-time applications due to its speed and interpretability via visual analysis. On the other hand, the LSTM model provides deeper insight into how a song evolves over time, which is crucial for rhythm- or beat-driven genre differentiation.

## 5.4 Summary

Both CNN and LSTM models are effective at classifying music genres, each excelling in different aspects. The CNN, using Mel spectrograms, is better at capturing timbral and harmonic textures. The LSTM, trained on MFCC sequences, shows strength in modeling rhythmic flow and temporal dependencies. Together, these models offer complementary perspectives and lay the groundwork for future hybrid architectures that leverage both spatial and temporal audio characteristics.

## 6. Conclusion

This study explored a hybrid deep learning model that integrates Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks for classifying music genres using the GTZAN dataset. By combining these two architectures, we aimed to take advantage of CNN's ability to extract spatial features from MFCC representations and LSTM's strength in modeling temporal patterns in audio signals.

The CNN component efficiently captured local time-frequency features, while the LSTM layers learned sequential relationships across time, which are essential for understanding musical structure. This combination led to more accurate genre classification than using either model individually.

Our approach achieved competitive performance, supported by the use of regularization techniques such as dropout and weight decay to prevent overfitting. These enhancements

improved the model's ability to generalize to new audio samples.

Overall, the CNN-LSTM model proved to be an effective and balanced solution for music genre classification. In future work, we plan to explore attention-based mechanisms, expand the model to real-time applications, and experiment with other datasets to validate its robustness.