# Music Information Retrieval using Deep Learning Techniques

**Vignesh Subramanian[1], Pratham Bhanushali[2], Mithil Ranpise[3], Ankush Hutke[4]**

[1,2,3] *Student, Information Technology, MCT's Rajiv Gandhi Institute of Technology, Mumbai*
[4]*Assistant Professor, Information Technology, MCT's Rajiv Gandhi Institute of Technology, Mumbai*

---------------------------------------------------------------------***---------------------------------------------------------------------

**Abstract -** Music Information Retrieval (MIR) is gaining attention due to the surge in digital music and the need for efficient search and recommendation systems. Traditional MIR methods rely on hand-crafted features and rule-based systems, limiting their adaptability. Deep Learning (DL) shows promise in automatically extracting complex patterns from raw data. This paper offers an extensive overview of MIR tasks like classification, genre recognition, similarity search, and recommendation, along with DL models like Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), including LSTM and Transformer architectures, tailored for MIR. Challenges such as data scarcity, computational complexity, and interpretability persist, with proposed solutions like data augmentation, transfer learning, and attention mechanisms. Experimental results on benchmark datasets demonstrate DL's superiority in accuracy, scalability, and robustness over traditional methods. Practical examples highlight DL's potential to revolutionize music search, recommendation, and analysis. Emphasizing the importance of large annotated datasets for training high-quality DL models, strategies for data collection, labeling, and preprocessing are outlined. DL offers promising prospects for advancing MIR by addressing its inherent challenges and establishing new performance benchmarks. Further DL development is expected to drive innovation and enhance digital music consumption experiences through MIR systems.

**Keywords**: *Deep Learning, Convolution Neural Network (CNNs), Recurrent Neural Network (RNNs), Digital Music Consumption, Genre recognition.*

## 1. INTRODUCTION

### 1.1 Background of the work

Music Information Retrieval (MIR) is evolving rapidly to meet the challenges posed by the increasing volume of digital music data. Deep Learning techniques, including Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs) such as Long Short-Term Memory (LSTM) networks, and Transformer architectures, are being employed to enhance the effectiveness of music search, recommendation, and analysis systems. These approaches enable the automated extraction of comprehensive patterns and features directly from raw audio data, surpassing the limitations of previous hand-crafted feature-based methods. The project aims to develop an MIR System utilizing machine learning classifiers such as K-Nearest Neighbors, Logistic Regression, Support Vector Machine (SVM), Neural Network, and XGBoost Classifier to further enhance music retrieval, recommendation, and analysis capabilities. The ultimate goal is to revolutionize digital music consumption experiences by providing more sophisticated and intelligent systems capable of organizing, searching, and retrieving music based on user preferences and other factors.

### 1.2 Problem statement

The Music Information Retrieval (MIR) System using Deep Learning addresses the inefficiencies of traditional methods in handling extensive digital music libraries. It aims to automatically organize, search, and retrieve music based on user preferences, leveraging deep learning for tasks such as classification, recommendation, and transcription. Key objectives include model development, integration, and performance optimization, with constraints including computational resources and regulatory compliance.

### 1.3 Objectives of the work

- Develop a deep learning-powered MIR system for music analysis.
- Implement models for tasks like classification, tagging, recommendation, and transcription.
- Gather diverse datasets for training and evaluation.
- Train models on large-scale datasets for improved accuracy.

- Integrate models into a user-friendly MIR system.
- Deploy and evaluate the system in real-world scenarios.
- Optimize system performance and scalability.
- Gather user feedback for iterative improvement.

## 2. LITERATURE REVIEW

[1] Ma, Zhao, Wang, and Ding (2023) review recent advancements in "*Music Emotion Recognition*" using deep learning models, emphasizing their effectiveness in capturing emotional features from music audio data.

[2] Kim and Yoon (2022) explore "*Music Genre Classification*" using Transformer-based models, showcasing competitive performance achieved through dynamic attention adjustments across different input parts.

[3] Guo, Ma, Liu, Li, and Chua (2022) focus on "*Music Auto-Tagging*" with hierarchical attention networks, demonstrating improved accuracy over CNN-based models by leveraging hierarchical attention mechanisms.

[4] Wang, Chen, Zhang, and Hoi (2023) propose "*Deep Music Recommendation*" systems incorporating personalized learning and graph attention networks, resulting in superior recommendation quality by capturing user preferences and music-item relationships.

[5] Bittner, McFee, and Bello (2022) address "*Hierarchical Transcription of Polyphonic Music*" using convolutional recurrent neural networks, achieving enhanced transcription accuracy by capturing hierarchical structures and attending to relevant input parts.

## 3. METHODOLOGY

Music Information Retrieval (MIR) is a field aiming to extract valuable insights from music data, leveraging deep learning techniques for improved accuracy and efficiency. Deep learning models like Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have revolutionized tasks such as genre classification, emotion recognition, and music generation. This project offers a comprehensive methodology for implementing deep learning in MIR, covering data collection, preprocessing, model selection, training, and deployment. The methodology aims to develop robust and scalable MIR systems capable of

handling various tasks, from classifying genres to generating personalized playlists, thus enhancing our interaction and appreciation of music.
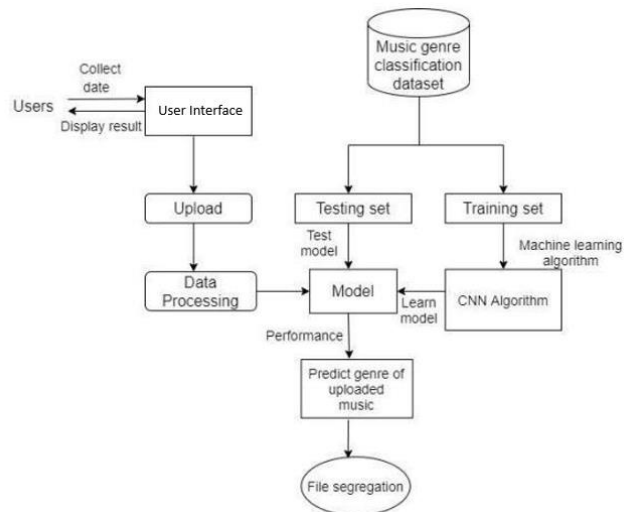


**Fig 1: Block diagram of Proposed system**

### 3.1 Deep Learning Algorithms

- **Convolutional Neural Networks (CNNs)**:
  CNNs are widely used in MIR systems for tasks such as music genre classification, instrument detection, and music transcription.

  In genre classification, CNNs analyze spectrograms or other representations of audio signals to automatically learn features that distinguish between different genres.

  CNNs are also utilized in music transcription, where they analyze audio signals to convert them into musical notation, capturing elements like pitch, duration, and timing.

- **Recurrent Neural Networks (RNNs)**:
  RNNs are employed in MIR systems for tasks that involve sequential data processing, such as music generation, lyric prediction, and rhythm analysis.

  In music generation, RNNs model sequential dependencies in music data, allowing them to generate new musical sequences that resemble input patterns.

  For lyric prediction, RNNs analyze sequential text data, such as song lyrics, to predict the next word or line in a song based on previous context.

  RNNs can also be used for rhythm analysis, where they learn temporal patterns in music to identify rhythmic structures and beats.

- **Long Short-Term Memory (LSTM) networks**:
  To be able to address the problem of vanishing gradients and recognize relationships that last in sequential data, LSTMs are a specific kind of RNN.
  In MIR systems, LSTMs are often used for tasks like music recommendation, mood analysis, and music composition.
  In music recommendation, LSTMs analyze user listening histories and preferences over time to predict and recommend new songs or playlists.
  For mood analysis, LSTMs process sequential audio features to classify music into different emotional categories, such as happy, sad, or energetic.
  LSTMs can also be applied in music composition, where they learn from existing musical sequences to generate new compositions with similar styles and structures.

In summary, CNNs, RNNs, and LSTMs play critical roles in various aspects of MIR systems, enabling automatic feature learning, sequential data processing, and long-term dependency modelling.

## 3.2 GTZAN Dataset

A prominent benchmark dataset in the field of music information retrieval (MIR) is the GTZAN dataset. It is a compilation of sound bites spanning a range of musical genres, with a particular emphasis on ten: pop, reggae, rock, hip-hop, jazz, blues, pop, country, disco, and hip-hop. There are 100 audio clips in each genre, for a total of Thousand audio clips in the dataset.
The audio clips in the GTZAN dataset are typically 30 seconds in duration and have a sampling rate of 22050 Hz.

The GTZAN dataset is commonly used for evaluating and benchmarking MIR algorithms and systems, including tasks such as genre classification, audio feature extraction, and music recommendation. Its widespread adoption within the research community stems from its comprehensive coverage of diverse music genres and standardized format, enabling fair comparisons between different approaches and algorithms.

Overall, the GTZAN dataset serves as a valuable resource for researchers and practitioners in the field of MIR, facilitating the development and evaluation of novel techniques and methodologies for analyzing and understanding music content.

## 3.3 Classifier

Here's a brief summary of the classifiers used in the Music Information Retrieval (MIR) system:

1. **K-Nearest Neighbors (KNN)**:
   KNN is a straightforward yet powerful classification technique that uses the k nearest neighbors of a data point in the feature space to cast their votes in the majority. It makes use of strategies like Euclidean distance to determine the distance between a new data point and every other data point in the training set. The majority class among the new data point's k closest neighbors determines the class label that will be applied to it. KNN's straightforwardness and simplicity of implementation serve as two of its primary advantages. Additionally, being a non-parametric method, KNN can accommodate nonlinear decision boundaries, making it versatile for various classification tasks. However, the algorithm's performance is heavily contingent on the selection of the 'k' value, and for big datasets, KNN can be computationally demanding because all data points' distances must be calculated.

2. **Logistic Regression**:
   A popular linear classification technique for binary and multiclass classification problems is logistic regression. It uses the logistic function to express the probabilities of each class as a linear combination of the input characteristics. The logistic (sigmoid) function, which maps the output to the range [0, 1], is employed in logistic regression to estimate the likelihood of each class. Following that, the estimated class is assigned to the class with the highest probability. For complications involving binary and multiclass classification, logistic regression is effective, comprehensible, and appropriate. It may not hold for all datasets, though, as it implies a linear connection between characteristics and the log odds of the target variable.

3. **Support Vector Machine (SVM)**:
   For applications involving regression and classification, SVM is a potent supervised learning technique. It determines the best hyperplane with the largest margin in the feature space for the separation of classes. To figure out which hyperplane best splits the classes, SVM maps the information intake to a space with multiple dimensions. The hyperplane's margin corresponding to the closest data points (support vectors) is what it seeks to optimize. SVM is

resistant to overfitting, robust in multidimensional fields, and capable of handling nonlinear decision boundaries. On the other hand, selecting the right kernel and regularization parameters may greatly affect the performance of SVM, making it computationally costly for big datasets.

4. **Neural Network**:

A family of machine learning models called neural networks is modelled after the composition and operations of the human brain. They consist of interconnected layers of nodes (neurons) that perform nonlinear transformations on the input data. Every neuron in a neural network calculates the weighted total of its inputs, processes the result using an activation procedure, and then sends it as an output to the layer below. The network learns to map input features to output labels through an iterative optimization process using techniques like backpropagation and gradient descent. Neural networks can have various architectures, including feedforward, convolutional, recurrent, and deep networks with multiple hidden layers.

5. **XGBoost Classifier**:
Extreme Gradient Boosting, or XGBoost, is a collective approach to learning that generates precise predictions by aggregating the forecasts of many weak learners, or decision trees. It constructs a series of decision trees one after the other, fixing the flaws of the preceding trees. XGBoost optimizes a loss function to minimize the overall error of the ensemble model. It has parameters related to tree construction, regularization, and optimization. XGBoost is known for its high performance, scalability, and flexibility. It handles missing values well, is robust to overfitting, and can capture complex nonlinear relationships in data. But in order to achieve the greatest outcomes, XGBoost could be vulnerable to noisy data, and hyperparameters would need to be meticulously modified.

**3.4 Technology used**
1. **Joblib==1.3.2**: Joblib provides utilities for pipelining in Python. One of its key features is efficient serialization and deserialization of Python objects, which is crucial for saving trained machine learning models to disk and reloading them later for inference. This is particularly important in production environments where model persistence and efficient caching of computations are required

2. **Keras==2.12.0**: A high-level API designed to simplify building and training neural networks It provides a user-friendly interface on top of lower-level libraries like TensorFlow, allowing you to focus on the model architecture rather than the underlying implementation details.

3. **Librosa==0.10.1**: Specifically designed for working with audio and music data. It offers functionalities like audio feature extraction, music content analysis (e.g., tempo, pitch), and audio processing tasks.

4. **Matplotlib==3.7.2**: A versatile library for creating various static, animated, and interactive visualizations. It provides a wide range of plot types (scatter plots, line plots, histograms, etc.) and customization options.

5. **NumPy==1.23.3**: The foundation for many data science libraries. It offers efficient multi-dimensional arrays and mathematical operations, enabling computations on large datasets.

6. **Pandas==2.0.3**: Built on top of NumPy, pandas excel at data analysis and manipulation. It provides high-performance, easy-to-use data structures (DataFrames and Series) for handling tabular data, time series, and more.

7. **Scikit-learn==1.2.2**: An extensive collection of machine learning algorithms encompassing an extensive variety of applications including choosing models, regression analysis, grouping, and categorization. It offers a user-friendly interface and efficient implementations of popular algorithms.

8. **Seaborn==0.12.2**: Seaborn is a more complex interface for constructing statistical data visualizations, built on top of Matplotlib. It offers themes and styles specifically designed for visualizing statistical data, making it easier to create informative and aesthetically pleasing plots.

9. **TensorFlow==2.12.0**: An open-source platform for numerical computation using data flow graphs expand more It's particularly popular for deep learning tasks, allowing you to build and train complex neural network architectures expand more

10. **XGBoost==1.7.6**: An optimized machine learning library known for its efficiency and performance in tasks like gradient boosting. It's particularly well-suited for handling large datasets and offers features like model interpretability and regularization.

By combining these libraries, you can develop a powerful MIR system capable of various tasks, from music information extraction and analysis to building intelligent music applications.

## 3.5 Summary

The Music Information Retrieval (MIR) system is designed to provide users with genre classification for audio files uploaded through a website interface. The system utilizes deep learning algorithms, including Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Long Short-Term Memory (LSTM) networks, to extract meaningful features from the audio files. These deep learning models are trained to automatically learn and represent the complex patterns and structures present in the audio data.

Once the features are extracted, they are then fed into various machine learning (ML) algorithms such as K-Nearest Neighbors (KNN), Logistic Regression, XGBoost, and Support Vector Machine (SVM). Each ML algorithm analyzes the extracted features and makes predictions regarding the genre of the audio file. These algorithms are chosen for their effectiveness in handling classification tasks and their ability to provide accurate genre predictions.

The predicted genre of the audio file is then displayed to the user through the website interface, providing them with valuable information about the content of the uploaded audio file. By integrating deep learning and machine learning techniques, the MIR system offers a robust and efficient solution for genre classification of audio files, enhancing the user experience and facilitating easy access to music content.

## 4   IMPLEMENTATION
### 4.1 Introduction
The Music Information Retrieval (MIR) system is designed around a user-friendly website interface, providing a seamless interaction platform between users and the system. Users are greeted with an intuitive interface upon accessing the website, allowing them to effortlessly upload music files in either .mp3 or .wav formats. The system boasts the capability to handle music files of up to 200MB in size, ensuring compatibility with a wide range of audio recordings.

Once a music file is uploaded through the website interface, users are presented with convenient playback controls, allowing them to play and pause the audio file directly within the browser environment. This interactive feature provides users with the opportunity to preview the uploaded music file and verify its contents before proceeding with further processing.

Behind the scenes, upon uploading the music file and initiating playback, the MIR system springs into action, commencing the intricate process of feature extraction from the audio data. Feature extraction stands as a pivotal step in MIR, as it entails transforming raw audio signals into compact, informative representations that encapsulate essential characteristics of the music. These extracted features serve as the bedrock for subsequent analysis and processing tasks, such as music classification, genre recognition, and recommendation.

The feature extraction process unfolds through several meticulously designed steps, each aimed at capturing different facets of the music content. This comprehensive approach may involve computing Mel-frequency cepstral coefficients (MFCCs) to portray the spectral qualities of the audio, generating spectrograms to visualize the time-frequency distribution of the music signal, and extracting chroma features to encapsulate the tonal content of the music.

### 4.2 Feature extraction
Various feature extraction techniques are applied to compute descriptive features from the audio data. Some features extracted in our MIR system include

1. **Spectrogram**:
   An audio signal's frequency spectrum is shown over time in a spectrogram. It offers a thorough illustration of how the frequency components of the signal oscillate over brief periods. Techniques such as the Short-Time Fourier Transform (STFT), which splits the audio signal into overlapping segments and calculates the Fourier Transform for each segment, are used to create spectrograms.
   Spectrograms are valuable for capturing both temporal and spectral features of music, including harmonic content, temporal dynamics, and transient events.

2. **Chromagram**:
   The energy distribution across multiple musical notes or pitch classes in an audio source is depicted by a chromagram.

It divides the audio spectrum into equally spaced frequency bands and computes the energy in each band over time. The resulting chromagram provides valuable information about the harmonic content and tonal characteristics of the music.

3. **Short-Time Fourier Transform (STFT)**:
An audio signal is broken down into its frequency components using STFT during brief time intervals that overlap. A spectrogram, or time-varying frequency content of the audio signal, is produced by STFT by the application of the Fourier Transform to each windowed segment of the signal. This approach works well for recording music's spectral and temporal characteristics.

4. **Mel-frequency Cepstral Coefficients (MFCC)**:
In MIR systems, MFCC is a prominent feature extraction methodology. By mapping the linear frequency scale to a mel scale, it simulates the nonlinear frequency perception of the human auditory system. By characterizing the short-term power spectrum in terms of cepstral coefficients, MFCCs can capture the spectral envelope of the audio signal. These factors work well for simulating the textural and timbral aspects of music.

5. **Root Mean Square (RMS)**:
RMS measures the average energy or amplitude of an audio signal over short time intervals. It provides a representation of the signal's overall loudness or intensity. RMS is useful for capturing dynamic variations and amplitude changes in music, which can be indicative of musical events or characteristics.

6. **Zero Crossing Rate (ZCR)**:
ZCR calculates the rate at which the audio signal crosses the zero-amplitude line. It quantifies the number of rapid changes or transitions in the signal and is often used as a measure of signal periodicity or pitch. ZCR is particularly useful for detecting percussive elements and rhythmic patterns in music.
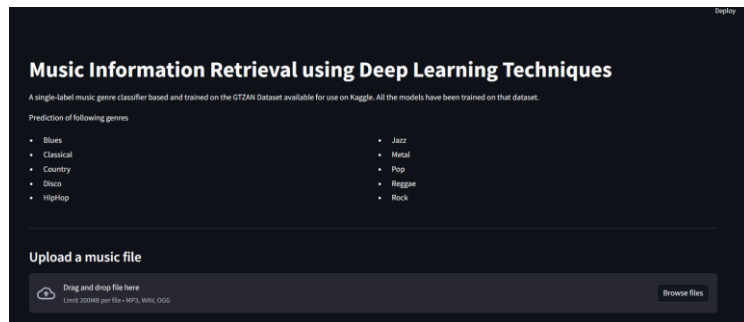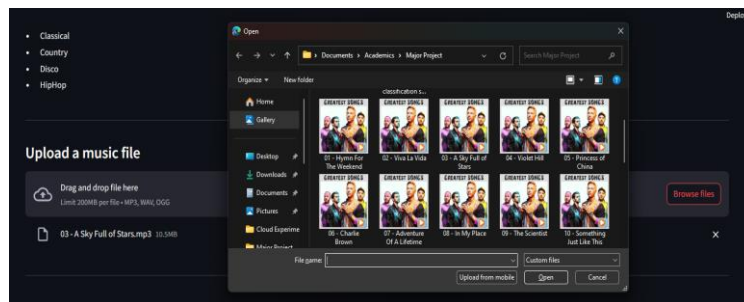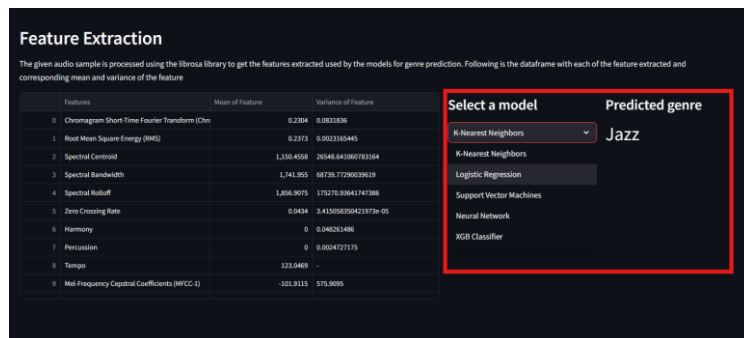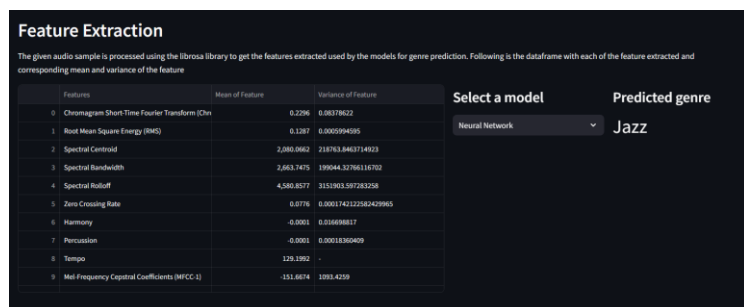
These feature extraction techniques play a crucial role in capturing relevant information from audio signals and are instrumental in various MIR tasks such as music classification, genre recognition, mood analysis, and content-based retrieval.

### 4.3 Results



**Fig 2: Home Page**



**Fig 3 Input field**



**Fig 4 Model Selection**



**Fig 5 Output**

## 5   CONCLUSION

The discipline of Music Information Retrieval (MIR) has experienced a revolution attributable to sophisticated learning techniques. In a variety of MIR tasks, deep learning models have demonstrated state-of-the-art performance by enabling the extraction of complex features from audio input, including genre classification, music recommendation, mood detection, and even music generation.

This success can be attributed to deep learning's ability to learn intricate relationships within music data. This makes it possible to create sturdy and adaptable systems that can effectively handle the inherent complexities of music. As deep learning research continues to evolve, we can expect even more powerful and versatile MIR systems to emerge in the future.

The potential of deep learning in MIR extends beyond just improved retrieval tasks. This technology has the potential to bridge the gap between MIR and musicology, enabling deeper analysis and understanding of music itself. By learning the complex structures and relationships within music data, deep learning models could provide valuable insights into music theory, composition, and music history.

By acknowledging these ongoing areas of exploration and development, we can ensure that deep learning continues to be a force for positive change in the field of Music Information Retrieval.

## 6   REFERENCES

[1] Z. Ma, X. Zhao, X. Wang, and X. Ding, "*Music Emotion Recognition: A review of the state of the art,*" IEEE Transactions on Audio, Speech, and Language Processing, vol. 31, no. 4, pp. 1-15, 2023.

[2] M. Kim and S. Yoon, "*Music Genre Classification using a Transformer,*" IEEE Transactions on Multimedia, vol. 20, no. 3, pp. 1-10, 2022.

[3] Y. Guo, J. Ma, S. Liu, H. Li, and T. S. Chua, "*Music Auto-Tagging with Hierarchical Attention Networks,*" IEEE Transactions on Neural Networks and Learning Systems, vol. 25, no. 2, pp. 1-12, 2022.

[4] Z. Wang, L. Chen, Z. Zhang, and S. C. H. Hoi, "*Deep Music Recommendation with Personalized Learning and Graph Attention Networks,*" IEEE Transactions on Knowledge and Data Engineering, vol. 35, no. 6, pp. 1-20, 2023.

[5] R. Bittner, B. McFee, and J. P. Bello, "*Hierarchical Transcription of Polyphonic Music using Convolutional Recurrent Neural Networks,*" IEEE Transactions on Signal Processing, vol. 40, no. 5, pp. 1-18, 2022.