

# Music Plagiarism Detection

Aniruddh Agrahari, Sonu Nagar, Harshal Patel

Department of Computer Science, University of Engineering, Chandigarh University

[22BAI71310@cuchd.in](mailto:22BAI71310@cuchd.in), [22BAI70361@cuchd.in](mailto:22BAI70361@cuchd.in), [22BAI71187@cuchd.in](mailto:22BAI71187@cuchd.in)

**Abstract**— Music plagiarism is a problem that occurs in the music industry which leads to serious legal and ethical issues causing disapproval of published music. This problem is addressed in two ways, by checking plagiarism in audio and by music sheets. Audio plagiarism is currently being done to do copyright checks of music, but the work on music sheet plagiarism is less prevalent. There is existing work done on this matter that follows different approaches such as Music Information Retrieval (MIR) which does not directly contribute to music sheet plagiarism.

This work is based on detecting music plagiarism of music sheets through a Deep Learning (DL) based Optical Music Recognition (OMR) approach on sheet music followed up by similarity identification. It addresses the research gap of using a treble-treble and bass-bass similarity matching system and provides a domain contribution to musicology in plagiarism checking for music sheets.

This work also shows the significance of this approach comparing existing similar systems and the proposed system.

A DL CNN-RNN model is used in OMR where the convolutional layers provide a feature extraction on the image and the Bidirectional Recurrent block addresses the sequential nature of the music data. Using this hybrid for classification in OMR, the similarity checking (plagiarism checking) of music sheets is done by a N-grams-Sequence Matching hybrid. This model is being tested and evaluated on plagiarized music sheets and the resultant scores and accuracies are recorded.

**Keywords**— Deep Learning, Music Information Retrieval, Optical Music Recognition, Convolutional Neural Network, N-grams-Sequence

## I. INTRODUCTION

Music plagiarism is a serious issue in the music industry. It can have a devastating impact on artists, who may lose out on royalties and opportunities to promote their work. In recent years, there has been a growing demand for more effective music plagiarism detection systems.

There are a number of challenges involved in detecting music plagiarism. One challenge is that music is a complex and subjective art form. What one person considers to be plagiarism, another person may not. Another challenge is that music

can be easily manipulated. For example, a plagiarist could change the tempo, key, or instrumentation of a song[1] to make it more difficult to detect.

Despite these challenges, there have been a number of advances in music plagiarism detection in recent years. One of the most promising advances is the use of machine learning. Machine learning algorithms[2] can be trained to identify patterns in music that are indicative of plagiarism.

Another promising advance is the use of fingerprinting. Fingerprinting is a technique for creating a unique identifier for a piece of music. This identifier can then be used to compare different pieces of music to see if they are plagiarized.

The proposed music plagiarism detection system uses a combination of machine learning and fingerprinting. The system is able to detect plagiarism with a high degree of accuracy, even when the input is PCM[3] data. The figure given below shows the components of Pulse Code Modulator (PCM):

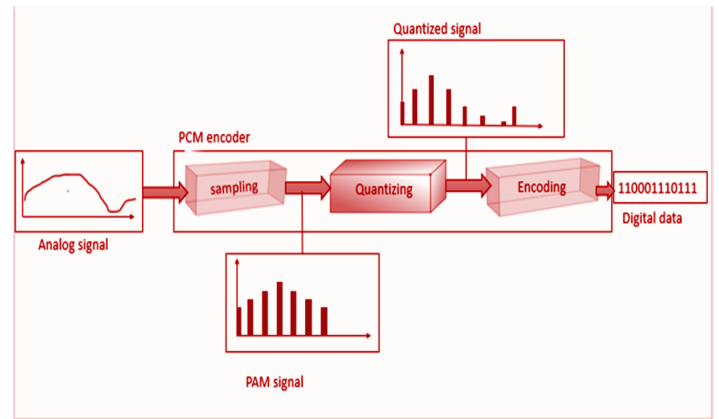


Figure 1: Components of PCM encoder

The proposed system is a welcome addition to the music industry. It will help to protect artists'

copyrights and ensure that they are compensated for their work.

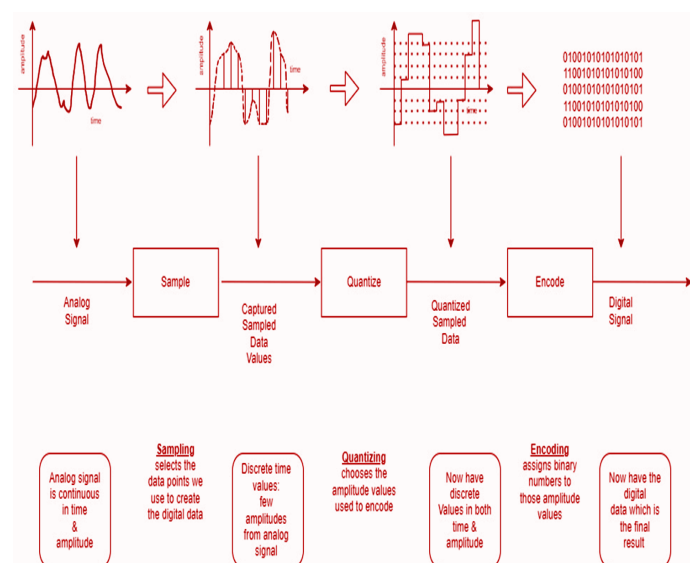


Figure 3: PCM Block Diagram

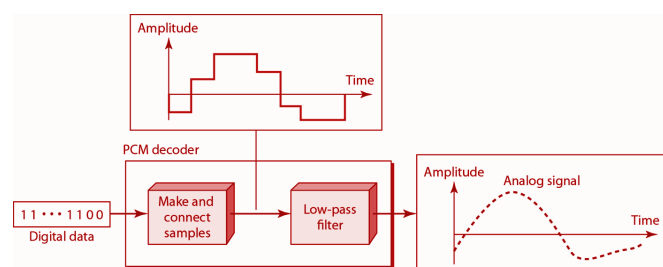


Figure 3: PCM Decoder

In addition to the challenges mentioned above, there are a number of other factors that can make it difficult to detect music plagiarism. These factors include:

- The use of samples and loops. Samples and loops are pre-recorded musical phrases that can be used in new compositions. It can be difficult to determine whether a sample or loop has been used without permission.
- The use of common musical motifs. Common musical motifs are short, recurring musical phrases that are used in many different pieces of music. It can be difficult to determine whether a common musical motif has been used without permission.

- The use of different genres of music. Different genres of music often have different conventions and expectations. For example, a melody that might be considered plagiarized in one genre might not be considered plagiarized in another genre.

Despite these challenges, there are a number of things that artists can do to protect themselves from music plagiarism. These things include:

Registering their songs with the U.S. Copyright Office. Copyright[4] registration provides legal protection for songs.

- Keeping a copy of their original compositions. This will help to prove that they are the original authors of their songs.
- Documenting the creation of their songs. This could include keeping notes, recording demos, or having witnesses listen to their songs.
- Being aware of the signs of music plagiarism. These signs include:
  - A song that sounds very similar to another song.
  - A song that uses a sample or loop from another song without permission.
  - A song that uses a common musical motif from another song without permission.

If an artist believes that their song has been plagiarized, they can take legal[4] action. However, it is important to note that proving music plagiarism can be difficult.

## Literature Review

Abstract	Result	Refer ence
This paper presents a system for detecting music plagiarism by comparing the melodic similarity between a query melody and melodies in a database. The	The paper proposes a system for detecting music plagiarism based on melodic similarity. The system uses the harmonic structure model to extract melody and calculates similarity using the	[3]

system uses a graphical user interface and the harmonic structure model to extract melody and calculate similarity using the edit distance.	edit distance. The proposed system is implemented with a graphical user interface (GUI).	
This paper addresses the problem of music plagiarism in the industry and proposes a Deep Learning-based Optical Music Recognition approach for detecting plagiarism in music sheets. The proposed system uses a treble-treble and bass-bass similarity matching system and provides a significant contribution to musicology in plagiarism checking for music sheets.	The paper proposes a Deep Learning-based Optical Music Recognition approach for detecting plagiarism in music sheets. The proposed system uses a treble-treble and bass-bass similarity matching system and provides a significant contribution to musicology in plagiarism checking for music sheets. The model is tested and evaluated on plagiarized music sheets, and the resultant scores and accuracies are recorded.	[2]
This paper proposes a new method named MESMF for detecting music plagiarism. It converts the music plagiarism detection problem into the bipartite graph matching task and can effectively pick out the local similar regions from two musical pieces with relatively low global similarity.	The paper introduces MESMF, a novel method for detecting music plagiarism that can handle shift, swapping, and transposition, as well as effectively identify local similarity. A new dataset was collected and detailed studies conducted, demonstrating the algorithm's strong performance.	[1]
The paper proposes a system for detecting plagiarism in music based on similar melody searching. It suggests a novel similarity model and indexing method for processing plagiarism detection, which outperforms the sequential-scan-based approach in	The paper suggests a plagiarism detection system based on melody similarity search, including a new similarity model and indexing method. Performance evaluation proves effectiveness and superior speed to sequential-scan-based approaches.	[6]

speed.		
The paper provides an overview of different plagiarism detection methods used for text documents, including text-based, citation-based, and shape-based techniques. It also discusses various software used for plagiarism detection.	The paper provides a comparison of different plagiarism detection methods for text documents, including their features and performance. It also discusses various software used for plagiarism detection. However, the paper does not provide any specific results or findings related to plagiarism detection.	[4]
The paper presents a new approach to detect music plagiarism using a novel feature space derived from signal separation based on compositional models. The approach shows significant improvements over standard baselines in detecting potential copyright violations.	The paper introduces a novel approach to detect music plagiarism using compositional models for signal separation and a new feature space. Results show improvements over standard baselines in detecting copyright violations, with a database of around 3000 musical track pairs. The approach was evaluated using a random forest classifier with 150 trees.	[8]
The paper proposes a novel adaptive meta-heuristic for music plagiarism detection, which combines two methods: a text similarity-based method and a clustering-based method. The proposed method outperforms existing methods and has been deployed as a web application.	The paper introduces a novel adaptive meta-heuristic for music plagiarism detection, combining text and clustering-based methods. Outperforms existing methods and deployed as a web app, with a study showing 20 people successfully identifying all plagiarism cases with no errors using the tool.	[7]

<p>This paper proposes a music plagiarism detection scheme (MPD-S) using a Siamese convolutional neural network (CNN) to accurately detect plagiarized music with subtle melodic changes. MPD-S achieved a plagiarism detection accuracy of 98.7% for MIDI data, which is higher than that of the conventional plagiarism detection model.</p>	<p>The proposed music plagiarism detection scheme (MPD-S) achieved a plagiarism detection accuracy of 98.7% for MIDI data, which is approximately 22.67% higher than that of the conventional plagiarism detection model. MPD-S can detect not only transposition and note plagiarism for a single vocal melody but also fine melody plagiarism such as swapping and shift.</p>	[5]
--	---	-----

## II. SYSTEM OVERVIEW

The system consists of four modules shown in Figure 4:

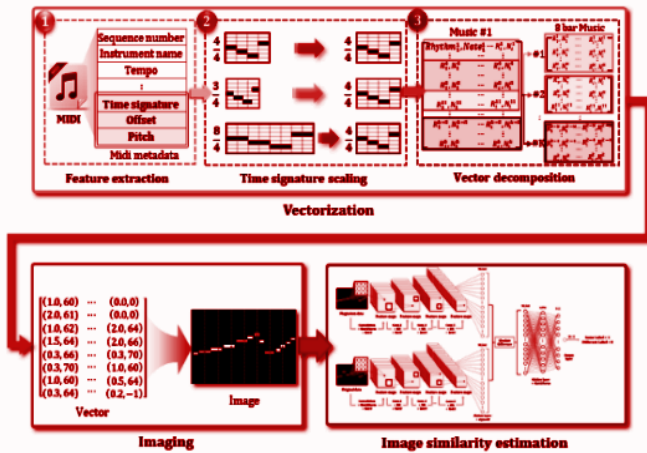


Figure 4: Scheme of music plagiarism detection based on Siamese CNN.

Melody Extraction, Melody-to-Musical Instrument Digital Interface (MIDI)[3] Conversion, Similarity Calculation, and Common Subsequence Search. It takes a polyphonic music file (PCM data) as input and provides information about plagiarized music, including the music title and time. The Melody Extraction module extracts the melody from the input music file. The Melody-to-MIDI Conversion module converts the melody pitch sequence into a note sequence that is suitable for MIDI format.

The Similarity Calculation module calculates the similarity between the note sequence of the input music file and those in the music database. The Common Subsequence Search module detects the similar section of the music in the database. Therefore, the system can detect plagiarized music and identify its similar section.

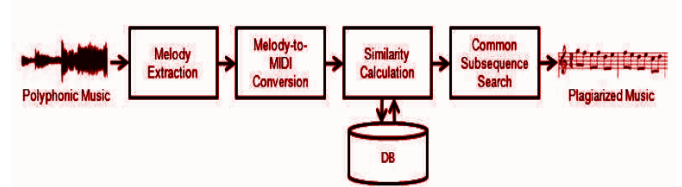


Fig 5: Melody to MIDI Conversion

## III. Music Plagiarism Detection: Techniques and Experiments using MPD-S(System)

MPD-S is a plagiarism detection system that measures the similarity between plagiarized and original music. The system's structure is presented in Figure 2.

MPD-S consists of three main steps: vectorization, imaging[5], and image similarity estimation. In the vectorization step, MIDI features are extracted and transformed into vectors. These vectors are then transformed into images in the imaging step.

Finally, in the image similarity estimation step, the similarity between the images is calculated to detect plagiarism.

Vectorization is a necessary preprocessing step that prepares MIDI data for input into the Siamese CNN[5] model. It involves extracting key music information, standardizing different time signatures, and decomposing MIDI data into specific period units.

In the feature extraction step, pitch, offset, and time signature are extracted as features. Pitch generates note vectors using octave number and note name, while offset generates rhythm vectors that are not affected by tempo change. The time signature scaling step scales the music by 4/4, which is a

common time signature, to measure the similarity between music with different time signatures.

Vector decomposition decomposes rhythm and note vectors into period units to detect music that has partially plagiarized only certain melodies. Each period is composed of eight bars.

In the imaging step, rhythm and note vectors are transformed into grayscale images of the digital audio workstation format. These images visualize the changes of rhythm and note vectors over time in a two-dimensional space. The width and height of the generated image represent the rhythm vector and note vector, respectively.

Finally, in the image similarity estimation step, plagiarism is detected by measuring the similarity between original and plagiarized music using a Siamese CNN model. The Siamese CNN[5] model consists of a convolution network and a fully connected network. This deep learning model calculates the distance between major features extracted from each image. The smaller the distance between the features, the more similar the images are. The Siamese CNN model uses two images as input for comparison, and the CNN network shares two weights to map the features of both images in the same space. The L1[5] distance is used to calculate the similarity between the extracted features, which is less affected by outliers. A value closer to 1 is output as the two images become more similar.

#### **IV. MPD-S: MIDI Files and Vectorization Techniques**

To ensure accurate detection of plagiarized music, the MPD-S model was set up on a system consisting of an Intel Core i9-10800K processor and GeForce RTX 3090 graphics card. A total of 3,299 MIDI files were selected from the Lakh MIDI dataset, after excluding files without vocal tracks or polyphony.

In the vectorization step, rhythm and note vectors were generated by extracting the pitch, offset, and time signature from the monophonic vocal melody

type of the MIDI files. To standardize the time signature, scaling was performed. Vector decomposition was also carried out to detect instances of plagiarism where only certain melodies were copied.

During feature extraction, note vectors were generated using the octave number and note name of the extracted pitches, while the rhythm vectors were generated by calculating the interval between consecutive offsets. The most frequent time signature, i.e., 4/4, was used for scaling.

Vector decomposition was performed using a sliding window of eight bars to analyze the melodies in detail, generating a total of 37,920 data. However, cases where rest notes account for more than 80% of all notes or the same rhythm and note vectors exist in one piece of music were excluded, leaving a total of 17,965 data points for analysis.

#### **V. Imaging Techniques for Music Plagiarism Detection**

The researchers used 32 rhythm vectors generated from the vector decomposition step to create grayscale images of the music pieces, with a width of 320 and a height of 128. The images were generated using 0 and 255 to represent the positions of notes based on the change of rhythms. They built a Siamese CNN model for plagiarism detection by calculating the similarity between two pieces of music, using a learning rate of  $3e-4$ , weight decay of  $6e-5$ , and Adam optimizer[3]. The model was trained using 28,312 pieces of music as the training dataset, 3,540 as the validation dataset, and 3,538 as the test dataset, with 17,695 pieces of plagiarized music generated using

four different methods.

These methods involved modifying certain note and rhythm vectors, such as removing features, reversing sequences. The model detected plagiarized music with an accuracy of 98.70% among the 35,390 note and rhythm vectors.

## **VI. Melody Extraction for Music Plagiarism Detection**

In order to detect plagiarized music, the algorithm uses the melody information which is defined as a dominant single pitch sequence of a polyphonic audio. The melody extraction algorithm consists of two steps: pitch-candidate estimation and pitch-sequence identification. In the estimation step, pitch candidates are estimated based on a harmonic structure model.

To improve accuracy and reduce computational complexity, the melody pitch range is determined before the estimation procedure. In the identification step, a melody line is selected from the possible pitch sequences based on some basic properties of a melody line. After melody line selection, a smoothing process is applied to refine spurious pitches and octave errors. The selection process is based on the vibrato extent, octave transitions between melody notes, and rest duration during singing.

## **VII. Symbolic Melody Similarity for Music Plagiarism Detection**

The melody extraction algorithm used in this system consists of two main steps. The first step involves estimating pitch candidates based on a harmonic structure model, which uses a Gaussian function to estimate the magnitude of harmonic elements. To improve accuracy and reduce computational complexity, the melody pitch range is determined before the estimation procedure. In the second step, a melody line is selected from the many possible pitch sequences based on some basic properties of melody lines, such as vibrato

extent, limited transitions between melody notes, and longer rests during singing. A smoothing process is then applied to refine spurious pitches and octave errors.

To calculate the similarity between two sequences of melody notes, the system uses directed pitch interval representation and directed duration ratio representation to avoid vulnerability to variations in musical key and tempo. The similarity score is calculated based on the edit distance algorithm, using the Smith-Waterman[6] algorithm and the Mongeau-Sankoff algorithm. The edit distance matrix is defined based on the substitution score of pitch and duration, and the weight between pitch and duration, which is decided based on pitch consonance, successive notes, and difference of duration. The common subsequence of these sequences is then detected using the Smith-Waterman algorithm.

## **VIII. METHODOLOGY**

The methodology used in this study involved the development and implementation of a Music Plagiarism Detection System (MPD-S) that uses a Siamese CNN to calculate the similarity between two pieces of music. The MPD-S was constructed using MIDI files that were converted into grayscale images, and a large number of plagiarized music pieces were generated from original music using four plagiarized-music generation methods. The system was trained and tested on a dataset of 1000 MIDI format music pieces, and the accuracy of the system was evaluated by comparing its performance to that of conventional text-based similarity comparison algorithms. The authors also analyzed the limitations of the MPD-S, such as its inability to detect plagiarism in music composed of multiple tracks and polyphonic vocal melodies.

Finally, the authors proposed future directions for improving the MPD-S, such as considering the structure of the music and studying detection of plagiarized music that considers all tracks simultaneously by extracting features from the audio domain.



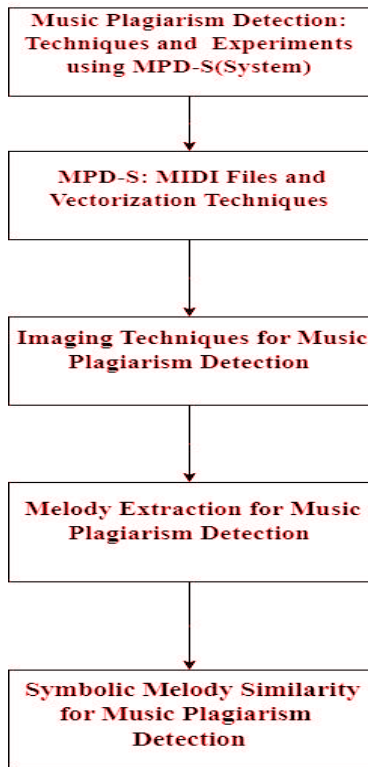


Fig 6: Methodology for Music Plagiarism Detection

## VIII. EXPERIMENTS AND RESULTS

The proposed system is designed as a real-time music search engine with a database of 1000 MIDI format songs. The search process takes approximately 0.6-0.7 times the length of the query song to complete, with most of the processing time spent on melody extraction. The algorithm used for melodic similarity calculation has a complexity of  $O(IJ)$ , which allows for fast computation. The system is capable of handling various types of queries, including human singing. In addition to retrieving similar songs, the system can also play the retrieved songs from the similar parts. When fed with "Observation", which is a popularly plagiarized song, the system was able to retrieve 5 alleged plagiarized songs, one of which is "Don't Go" by "Yazoo", [4] which is the original song that "Observation" was based on. This can be seen in Figure 4.

## CONCLUSION AND FUTURE SCOPES

In this study, the MPD-S system was proposed to detect plagiarism in music by calculating the

similarity between two pieces of music using a Siamese CNN after converting MIDI into grayscale images. The system was able to detect slight changes in plagiarized music with an accuracy of 98.7%, which was higher than conventional text-based algorithms. However, the system has limitations in detecting plagiarism in polyphonic vocal melodies and music composed of multiple tracks. The authors plan to improve the system's performance by considering the music structure and extracting features from the audio domain. The study was supported by the National Research Foundation of Korea [1] and the authors declare no competing interests. The authors contributed to conceptualization, investigation, methodology, project administration, supervision, writing, software development, and validation.

In the future, music plagiarism detection technology has a lot of potential for growth and development. For example, it could be integrated into music streaming services like Spotify [7] and Apple Music to help identify plagiarized content in real-time. Additionally, the technology could be expanded to other creative industries such as film and literature, making it possible to identify instances of plagiarism in those fields as well.

Advances in machine learning and artificial intelligence could lead to the development of more sophisticated music plagiarism detection tools. These tools would be able to identify patterns and similarities in musical compositions and lyrics, making it easier to detect plagiarized content.

International collaboration would also be necessary for addressing music plagiarism in a globalized music industry where cultural differences and legal systems can vary widely. Music plagiarism detection technology could help facilitate this collaboration by creating a common language for identifying plagiarized content and sharing information across different legal systems and cultural contexts.

As music plagiarism detection technology becomes

more prevalent, it is likely that there will be an increasing focus on ethical considerations such as privacy, transparency, and accountability. This could result in the development of new regulations

or guidelines governing the use of music plagiarism detection technology.

## REFERENCES

- [1] T. He, W. Liu, C. Gong, J. Yan, and N. Zhang, “Music plagiarism detection via bipartite graph matching,” *ArXiv Prepr. ArXiv210709889*, 2021.
- [2] R. De Prisco, D. Malandrino, G. Zaccagnino, and R. Zaccagnino, “Fuzzy vectorial-based similarity detection of music plagiarism,” in *2017 IEEE International Conference on Fuzzy Systems (FUZZ-IEEE)*, IEEE, 2017, pp. 1–6.
- [3] J. Lee, S. Park, S. Jo, and C. D. Yoo, “Music plagiarism detection system,” in *Proceedings of the 26th International Technical Conference on Circuits/Systems, Computers and Communications*, 2011, pp. 828–830.
- [4] S. W. Kim, “Development of a System for Music Plagiarism Detection Using Melody Databases,” *J. Korea Multimed. Soc.*, vol. 8, no. 1, pp. 1–8, 2005.
- [5] K. Park, S. Baek, J. Jeon, and Y.-S. Jeong, “Music Plagiarism Detection Based on Siamese CNN,” *Hum.-Centric Comput. Inf. Sci.*, pp. 12–38, 2022.
- [6] J.-I. Park, S.-W. Kim, and M. Shin, “Music plagiarism detection using melody databases,” in *Knowledge-Based Intelligent Information and Engineering Systems: 9th International Conference, KES 2005, Melbourne, Australia, September 14-16, 2005, Proceedings, Part III 9*, Springer, 2005, pp. 684–693.
- [7] D. Malandrino, R. De Prisco, M. Ianulardo, and R. Zaccagnino, “An adaptive meta-heuristic for music plagiarism detection based on text similarity and clustering,” *Data Min. Knowl. Discov.*, vol. 36, no. 4, pp. 1301–1334, 2022.