

Musical Genre Classification

Joey Zheng

joeyz@mit.edu

Daniel Kim

dyk0518@mit.edu

A. Abstract

Musical Genre Classification (MGC) is the process that predicts a musical genre given a musical audio input. Previous work has used both visual and audio features representation of the input, achieving as high as 93.4% accuracy on the GZTAN dataset. When we limit the scope of feature extraction to only visual representation, the best accuracy achieved by the previously published models is 89.30%. In this paper, we show that higher accuracy can be achieved by extracting more and combining multiple audio-visual features, as well as making modifications to the model architecture. By using a combined audio-visual feature of Mel-spectrogram, MFCCs, and Chromagram with a hybrid CNN + Bi-GRU architecture, we achieved an empirical accuracy of 89.94%. This accuracy is higher than the accuracy achieved by other published models that conduct MGC on the GTZAN dataset using audio-visual features only.

B. Introduction

Think of your favorite genre of music. Have you ever wondered, what it is about that genre that you specifically enjoy? Whether it be improvisation in jazz or simple harmonies in country, each genre has specific musical features that attract its audience. With recent signal processing and machine learning development, we can computationally extract and identify these key audio features from musical pieces [11].

In this paper, we aim to use these extracted features from musical audio and use them to predict the musical genre. This process is known as Musical Genre Classification (MGC). With countless musical pieces produced every day in the prolific music industry, automating the classification and analysis tasks can benefit the field as a whole. There are many applications for this concept. Digital music services, such as Spotify, use MGC in the process of music categorization and recommendation, as well as providing data sources that could be analyzed through MGC in third-party research [3]. Furthermore, MGC helps expand and develop the broader concept of Music Information Retrieval (MIR), a multidisciplinary field that focuses on the extraction, analysis, organization, and retrieval of music-related information [4]. In this paper, we explore the current state of MGC using visual features, and how we improved

it with data augmentation and architecture modifications of the model.

C. Background and Related Work

The standard dataset that is used for MGC is the GZTAN dataset, which has become the benchmark for musical analysis [13]. Hence, we also conduct our study using the GZTAN dataset. Currently, one of the best-performing MGC models for this dataset was introduced by Dai et al., which uses Mel-Frequency Cepstral Coefficients (MFCCs) and other audio features as the inputs with Deep Neural Networks to achieve 93.4% accuracy [2]. Dai et al. use two separate pipelines, Visual Feature Extraction (VFE) and Audio Feature Extraction (AFE) modules, during its data feature extraction process.

In our paper, we narrow down the scope to the MGC models that focus on using the image-level feature extraction of the musical audio input. By excluding the audio feature extraction in the process, we limit our study to how well computer vision techniques can be used for MGC. Within this narrowed scope, the best-performing model has been published by Ashraf et al. [1]. In their study, Ashraf et al. compare the performance based on two visual features from the audio: Mel-spectrogram and MFCCs, as well as a hybrid architecture of CNN and variants of RNN such as LSTM, Bi-LSTM, GRU, and Bi-GRU. Empirically, the best accuracy of 89.30% was achieved through the proposed hybrid architecture of CNN and Bi-GRU using Mel-spectrogram.

There are, however, a few ways that we suggest could improve this MGC using audio-visual features. We highlight two main ways this improvement could be achieved, data augmentation and architecture modification, which are also the main contributions of our paper:

- In addition to two audio-visual features used in the previous state-of-the-art model, Mel-spectrogram and MFCCs, we also use Chromagrams. Moreover, rather than training the model using each audio-visual feature individually, we also try training the model with combined audio-visual features. This way, each audio sample is represented with more diverse, informative features in the model.
- Compared to the previous state-of-the-art model, we reduce the number of CNN layer blocks, which helps with the computational efficiency of the model. We also replace the original RNN layer, which consisted of GRU, Bi-GRU, and GRU, with two Bi-GRU layers that led to better performance of the model.

With these changes, our model achieved the best accuracy of 89.94%, surpassing the performance of the previous

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

108 MGC models that use audio-visual features. In the follow-
109 ing sections, we dive deeper into our methodology.

111 D. Methodology

113 D.1. Hypothesis

114 Our hypothesis is that a musical genre classification
115 model that trains using combined audio-visual features will
116 outperform those that train using only one audio-visual fea-
117 ture. We base our hypothesis on two key observations.
118 Firstly, once audio-visual features are extracted, they can
119 be simply treated as images. This implies that image classi-
120 fication techniques, such as CNNs, can be applied to solve
121 this classification problem. Secondly, we observe that audio
122 is multi-dimensional. One could extract many features
123 from a piece of audio (especially musical audio) such as
124 frequency, timbre, rhythm, dynamic, etc. Previous work in
125 musical genre classification tends to focus on one audio-
126 visual feature at a time. We want to investigate the poten-
127 tial of combining audio-visual features for training musical
128 genre classification models.

130 D.2. Dataset

131 Our model was trained on the GZTAN dataset. This is
132 a publicly available dataset consisting of 10 labeled gen-
133 res: blues, classical, country, disco, hiphop, jazz, metal,
134 pop, reggae, and rock. Each genre comes with 100 audio
135 tracks each 30 seconds in length. The audio tracks are
136 all 22050Hz mono 16-bit audio files in .wav format with
137 a storage size of about 1.3MB. These audio files were col-
138 lected in 2000-2001 from various recording conditions such
139 as personal CDs, radio, microphone recordings, etc. Addi-
140 tionally, the dataset includes 2 CSV files with statistics on
141 the full audio file (30 seconds) and split audio file (3 sec-
142 onds). Among some of the statistics are spectral centroid
143 means, spectral centroid variances, MFCC means, MFCC
144 variances, tempo, and more. However, we did not utilize
145 these statistics, which we will discuss in the future work
146 section. We chose the GZTAN dataset as it serves as the
147 musical analysis benchmark in the field of musical genre
148 classification.

150 D.3. Audio-Visual Features

151 Audio-visual features refer to features extracted from the
152 audio that has visual components and can be analyzed with
153 computer vision techniques. More specifically, this paper
154 works with features such as spectrograms, which span the
155 time axis and frequency axis with the frequency intensity at
156 each point. This is in contrast with features such as spec-
157 tral centroids, which measure the weighted mean of the fre-
158 quencies at a given time.

160 This paper utilizes three audio-visual features: Mel-
161 spectrograms, MFCCs, and Chromagrams.

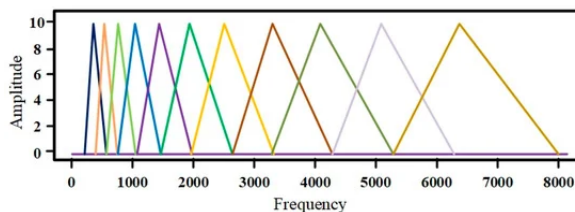
162 D.3.1 Mel-spectrogram

163 Spectrograms have high utility for visualizing audio. Audio
164 can be thought of as a combination of varying amplitudes of
165 frequency over time. A spectrogram decomposes audio into
166 its time and frequency components. Specifically, it maps a
167 given time and frequency to the intensity of that frequency
168 at that time. This is achieved by applying the Short-Time
169 Fourier Transform (STFT). The STFT is a Fourier transform
170 performed on smaller windows, or segments, of the audio.
171 This allows for the extraction of localized frequency con-
172 tent, which is more suitable for the analysis of frequently
173 varying audio such as music.

174 Mel-spectrogram is a transformed spectrogram using the
175 mel-scale, a non-linear scale that better approximates the
176 perception of the human auditory system. The approximate
177 formula for the mel-frequency of linear frequency in hertz,
178 f , is as follows:

$$180 \text{mel}(f) = 2595 * \log_{10} \left(1 + \frac{f}{700} \right)$$

181 To construct a Mel-spectrogram, a parameter for the
182 number of mel-bands, n_{mel} , is chosen (in our models, we
183 used $n_{mel} = 16$). A higher number of mel-bands implies a
184 more detailed representation of the audio data (though too
185 much may lead to overfitting). Using the mel-scale and
186 n_{mel} , a mel-bank filter is constructed and then applied to
187 the audio data to generate the Mel-spectrogram (Figure 1).



191 Figure 1. Example mel-bank filter with $n_{mel} = 11$

192 An example of a Mel-spectrogram on one of our audio
193 files is shown in Figure 2.

201 D.3.2 Mel Frequency Cepstral Coefficient (MFCC)

202 The MFCC is a compact and low-dimension representation
203 of the audio data by applying the Discrete Cosine Trans-
204 form (DCT) to the Mel-spectrogram. As a result, the most
205 significant cepstral coefficients are extracted. Cepstral coef-
206 ficients represent the spectral envelope of the audio data and
207 were found to be beneficial as features for machine learning
208 models. To generate the MFCC, a parameter for the number
209 of cepstral coefficients, n_{mfcc} , is chosen (in our models, we
210 used $n_{mfcc} = 13$). An example of an MFCC on one of our
211 audio files is shown in Figure 3.

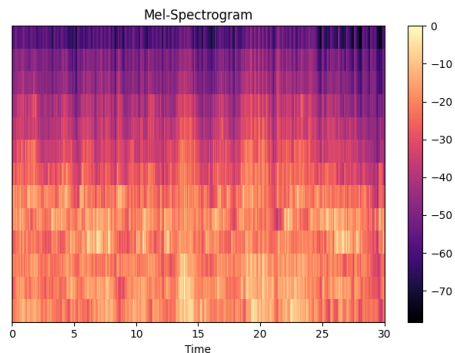


Figure 2. Mel-spectrogram of an audio labeled as classical

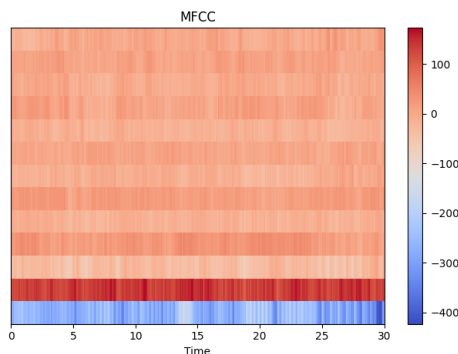


Figure 3. MFCC of the same audio labeled as classical

D.3.3 Chromagram

Chromagram is a type of spectrogram that represents the energy distribution of the 12 standard pitch classes of modern Western music by applying the STFT with respect to the 12 chroma values. Consequently, Chromagrams are invariant to octave differences, timbre, instrumentation, etc.

An example of an Chromagram on one of our audio files is shown in Figure 4.

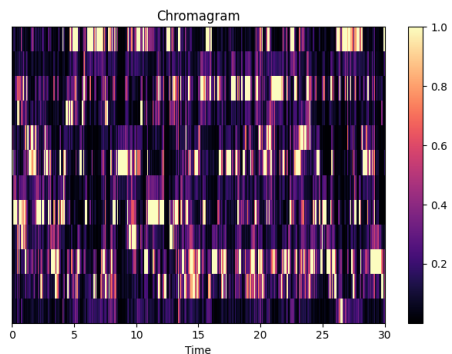


Figure 4. Chromagram of the same audio labeled as classical

D.4. Method

D.4.1 Preprocessing

To prepare the data for training, we performed data augmentation on the dataset. First, each 30-second audio file aforementioned was further split into 5 segments of 6 seconds each. This augmentation increases the amount of labeled data by a factor of 5. Six-second audio files are long enough such that the overarching context (the genre in our case) is maintained while short enough to allow for a higher prediction rate. As a result, there will be 5000 independent audio data each having a genre label and 6 seconds in length.

After the segmentation, audio-visual features are extracted from the audio data. Each 6-second audio file will generate a Mel-spectrogram, MFCC, and Chromagram. These audio-visual features, along with its genre label will be stored in a JSON file to be used later. When we use multiple of these extracted features to represent an audio input, we combined (through concatenation) those features before inputting them into the model. This process of extracting audio-visual features from input audio is illustrated in Figure 5.

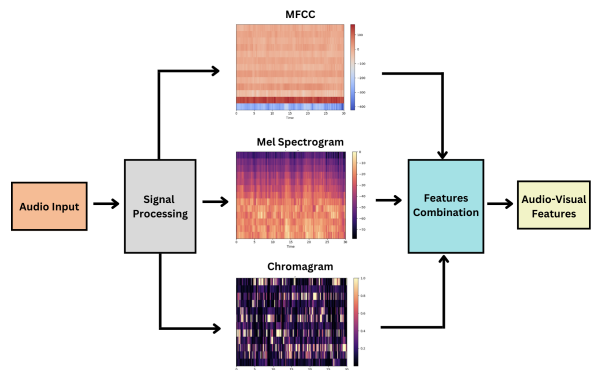


Figure 5. Audio input data processing diagram

D.4.2 Model

We investigated the performance of our model based on 5 different model architectures and 5 different sets of audio-visual features. Specifically, the 5 architectures are CNN, LSTM, Bi-GRU, CNN + LSTM, and CNN + Bi-GRU. The 5 audio-visual features are MFCC, Mel-spectrogram, Chromagram, MFCC + Mel-spectrogram, and MFCC + Mel-spectrogram + Chromagram. We trained and analyzed the performance of each combination of model architecture and audio-visual features, leading to 25 different models that were tested in total.

Below in Figures 6 7 8, we display the architectures of CNN, CNN + LSTM, and CNN + Bi-GRU respectively (the

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

architectures for our Bi-GRU and LSTM models individually are not displayed as they use the same parameters as in the hybrid architectures).

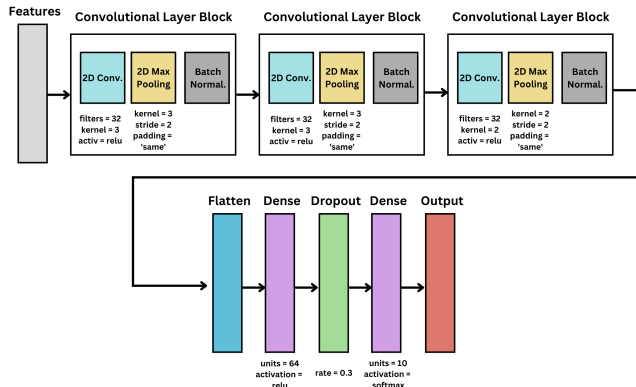


Figure 6. Proposed CNN architecture

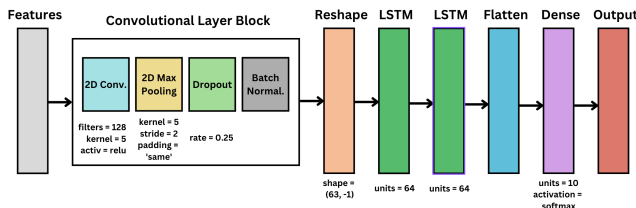


Figure 7. Proposed CNN + LSTM hybrid architecture

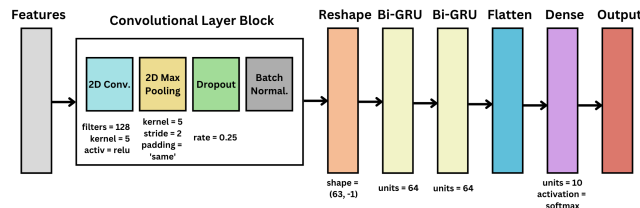


Figure 8. Proposed CNN + Bi-GRU hybrid architecture

D.4.3 Training

To prepare for model training, the audio-visual features JSON file is loaded and then partitioned into a training set (60%), validation set (15%), and testing set (25%). The labels are then one-hot encoded into 10 perpendicular vectors. In our models, we used the Adam optimizer with a learning rate of 0.0001 as an optimizer and categorical cross-entropy for the loss function. As hyperparameters, each model used a batch size of 32 and trained for 30 epochs. In our largest model (CNN + Bi-GRU trained on MFCC + Mel + Chroma features), each epoch to approximately 230 seconds to complete.

D.5. Performance Metrics

The performance of our model is measured by the cumulative test accuracy of the ground truth genre label of the audio data versus the predicted genre label of the same audio data. In other words, the metric is the accuracy of the model, namely the proportion of correctly classified audio data.

E. Experimental Results and Discussion

The summary of the accuracy in our study is presented in Table 1. Moreover, the graph that compares the performance of different audio-visual features on the best-performing architecture, CNN + Bi-GRU hybrid, is shown in Figure 9.

		Audio-Visual Features				
		MFCC	Mel	Chroma	MFCC + Mel	MFCC + Mel + Chroma
Model Architecture	CNN	81.57%	77.70%	55.00%	84.78%	82.47%
	LSTM	81.49%	66.91%	40.17%	74.10%	75.45%
	Bi-GRU	80.34%	75.22%	38.71%	79.83%	80.28%
	CNN + LSTM	84.72%	81.97%	62.81%	87.30%	86.40%
	CNN + Bi-GRU	89.89%	84.33%	68.37%	88.54%	89.94%

Table 1. Summary of the accuracy. For each architecture (row), the highest accuracy achieved is bolded. The overall highest accuracy is in blue.

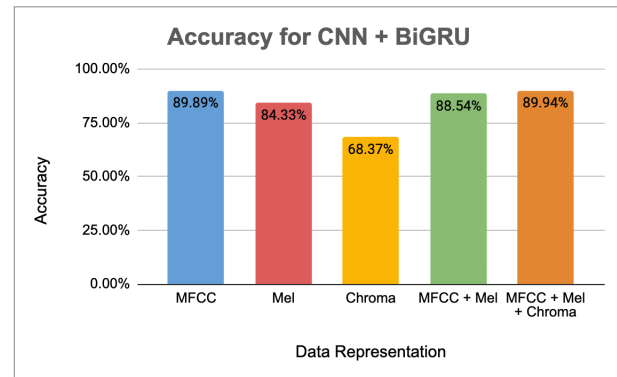


Figure 9. Accuracy with CNN + Bi-GRU hybrid Architecture

The findings indicate that having audio represented in multiple audio-visual features generally leads to better performance. This can be explained by the fact that with more audio features, the model has more characteristics of each genre of music to learn from. There were some

architectures, however, where this causal relationship between more features and better performance was not observed. For example, under LSTM architecture, the model achieved 81.49% accuracy with just the MFCC feature, higher than 74.10% (MFCC + Mel) or 75.45% (MFCC + Mel + Chroma) achieved with multiple features. Although we would like to explore this phenomenon more in-depth in the future, we hypothesize that extra features that LSTM fails to learn the relevance of have led to an over-fitted model.

In terms of model architecture, we see that the hybrid of CNN and Bi-GRU layers perform the best in every feature representation. This hybrid that uses both CNN and RNN layers perform better than using each of them alone for it extracts both spatial and temporal information from the audio-visual features. Specifically, with the feature representations of MFCC, Mel, and Chroma, the model achieved 89.94% accuracy in the test set. The plot for Accuracy and Error against Epoch in Figure 10 does not show signs of overfitting since the testing accuracy does not deviate too much from the training accuracy with the increasing number of epochs. This means that the weights generalize well to the dataset. This accuracy is higher than the accuracy achieved by other published models that conduct MCG on the GTZAN dataset using audio-visual features only. The comparison of our model's performance and other state-of-the-art models is shown in Table 2.

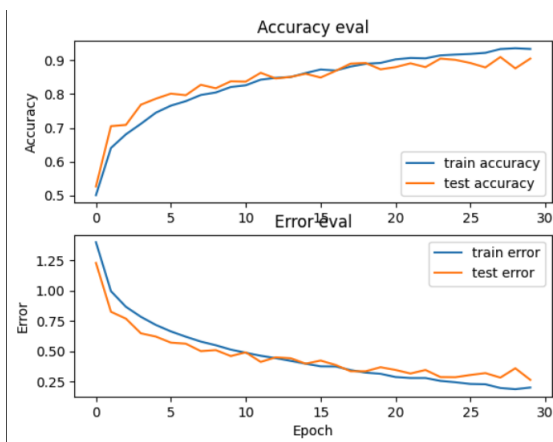


Figure 10. Accuracy and Error of CNN + Bi-GRU hybrid architecture with MFCC + Mel + Chroma features

F. Conclusion and Future Work

In this paper, we found that combining audio-visual features for training was beneficial for musical genre classification, which aligns with our hypothesis. More specifically, our CNN + Bi-GRU architecture achieved an accuracy of 89.94% using a combination of MFCC + Mel + Chroma features on the GTZAN dataset. Our model outperformed

Table 2. Comparison of State-of-the-Art Models

Method	Accuracy
George Tzanetakis [13]	61.00%
G. Sun et al. [12]	66.40%
A Heaki et al. [7]	70.60%
Nilesh M. [10]	77.78%
Praseneet Fulzele et al. [6]	89.00%
N. Farajzadeh [5]	86.00%
Pradeep Kumar D. et al. [9]	86.00%
Jan Jakubik [8]	87.70%
Mohsin Ashraf et al. [1]	89.30%
Our Model	89.94%

other state-of-the-art models that only used audio-visual features such as the 89.30% accuracy achieved in January 2023 [1].

Due to time and computing constraints, we understand there are many potential improvements that could be made to our model. For future works, we will incorporate other audio features such as spectral centroid, tempo extraction, percussive features, harmonic feature, etc. as presented in Jinliang L. et al 2021 [2], which achieved an accuracy of 93.4% by using both Visual Feature Extraction module and Audio Feature Extraction module. One approach is to utilize the statistics provided in the GTZAN dataset. We will also try higher resolution of audio-visual feature extraction by increasing the number of mel-bands for generating Mel-spectrograms (this was again due to limited computational power). Additionally, we will tweak and optimize hyperparameters (i.e. number of layers, epoch, etc.), train on other audio datasets (such as the Million Song Dataset), and investigate other ML architecture models (e.g. transformers).

G. Individual Contributions

- Literature reviews (Joey Zheng and Daniel Kim)
- Audio data download & feature extraction & preprocessing (Joey Zheng)
- Model layer architecture design & implementation (Daniel Kim)
- Training and testing of different models (Joey Zheng and Daniel Kim)
- Final report & presentation (Joey Zheng and Daniel Kim)

References

- [1] M. Ashraf, F. Abid, I.U. Din, J. Rasheed, M. Yesiltepe, S.F. Yeo, and M.T. Ersoy. A hybrid cnn and rnn variant model for music classification. *Applied Sciences*, 13, 2023. 1, 5
- [2] Jia Dai, Wen-Ju Liu, Chongjia Ni, Like Dong, and Hong Yang. "multilingual" deep neural network for music genre classification. 09 2015. 1, 5

540	[3] Tyler Dammann and Kevin Haugh. Genre classification of	594
541	spotify songs using lyrics, audio previews, and album art-	595
542	work, 2017. 1	596
543	[4] J Stephen Downie. Music information retrieval. <i>Annual re-</i>	597
544	<i>view of information science and technology</i> , 37(1):295–340,	598
545	2003. 1	599
546	[5] N. Farajzadeh, N. Sadeghzadeh, and M. Hashemzadeh. Pmg-	600
547	net: Persian music genre classification using deep neural net-	601
548	works. <i>Entertainment Computing</i> , 44:100518, 2023. 5	602
549	[6] P. Fulzele, R. Singh, N. Kaushik, and K. Pandey. A hybrid	603
550	model for music genre classification using lstm and svm. In	604
551	<i>Proceedings of the 2018 Eleventh International Conference</i>	605
552	<i>on Contemporary Computing (IC3)</i> , pages 1–3. IEEE, 2018.	606
553	5	607
554	[7] A. Heakl, A. Abdelgawad, and V. Parque. A study on broad-	608
555	cast networks for music genre classification. In <i>Proceedings</i>	609
556	<i>of the 2022 International Joint Conference on Neural Net-</i>	610
557	<i>works (IJCNN)</i> , pages 1–8, Padua, Italy, July 2022. 5	611
558	[8] J. Jakubik. Evaluation of gated recurrent neural networks	612
559	in music classification tasks. In <i>Proceedings of the 38th In-</i>	613
560	<i>ternational Conference on Information Systems Architecture</i>	614
561	<i>and Technology—ISAT 2017</i> , pages 27–37. Springer, 2017.	615
562	5	616
563	[9] D.P. Kumar, B.J. Sowmya, Chetan, and K.G. Srinivasa. A	617
564	comparative study of classifiers for music genre classifica-	618
565	tion based on feature extractors. In <i>2016 IEEE Distributed</i>	619
566	<i>Computing, VLSI, Electrical Circuits and Robotics (DIS-</i>	620
567	<i>COVER)</i> , pages 190–194. IEEE, 2016. 5	621
568	[10] N. M. Patil and M. U. Nemade. Music genre classification	622
569	using mfcc, k-nn and svm classifier. <i>Int. J. Comput. Eng.</i>	623
570	<i>Res. Trends</i> , 4:2349–7084, 2017. 5	624
571	[11] Garima Sharma, Kartikeyan Umapathy, and Sridhar Krish-	625
572	nan. Trends in audio signal feature extraction methods. <i>Ap-</i>	626
573	<i>plied Acoustics</i> , 158:107020, 2020. 1	627
574	[12] G. Sun. Research on architecture for long-tailed genre com-	628
575	puter intelligent classification with music information re-	629
576	trieval and deep learning. <i>Journal of Physics: Conference</i>	630
577	<i>Series</i> , 2033:012008, 2021. 5	631
578	[13] G. Tzanetakis and P. Cook. Musical genre classification of	632
579	audio signals. <i>IEEE Transactions on Speech and Audio Pro-</i>	633
580	<i>cessing</i> , 10:293–302, 2002. 1, 5	634
581		635
582		636
583		637
584		638
585		639
586		640
587		641
588		642
589		643
590		644
591		645
592		646
593		647