

---

# Deep Learning Methods Explored for Music Genres Classification

---

Haopeng Wu  
University College Dublin  
haopeng.wu.ucdconnect.ie

## Abstract

Music plays an important part in people's daily life. The classification of music genres is important for music recommendation. **Deep learning** methods have shown promising results in the automatic **music genres classification**. In this document, various deep learning techniques for music genre classification are explored and compared, including Multilayer Perceptrons (**MLPs**), Convolutional Neural Networks (**CNNs**), Recurrent Neural Networks (**RNNs**), and **hybrid architectures**. The impact of different input representations are also investigated, such as **Mel Spectrograms** and **MFCCs**, on the performance of these models. **Transferred learning** is employed to use pretrained models.

A lot of interesting findings are found along the way of the exploring journey. The most important finding is how easy it is for one to build a decent deep learning model to address a problem in a particular field with the help of the abundant resources online.

## 1 Introduction

Music is a universal language that has been enjoyed by people across the globe for centuries. Over time, music has evolved and diversified, resulting in the emergence of various music genres that cater to different tastes and preferences. While humans have been able to classify music into various genres based on their unique characteristics, the advent of deep learning technology has made it possible for machines to do the same.

Deep learning is a branch of machine learning that uses neural networks to learn and make decisions based on input data. With the help of deep learning algorithms, it is possible to classify music genres automatically by different music features.

Music genre classification by deep learning has various applications, including music recommendation systems, audio-based search engines, and music analysis tools. In recent years, researchers have explored various deep learning techniques such as MLPs, CNNs, and RNNs to classify music genres with high accuracy by different music data.

This document explores various deep learning techniques to address music genres classification and the differences are compared and analysed. Specifically, 5 deep learning models are implemented based on 3 different audio features. Each model is comprehensively hypertuned against architecture, activation function, units in layers, dropout rate, learning rates, batch size, and epochs. The 5 models are MLP, VGG16, YAMNet, LSTM, and a hybrid of YAMNet and LSTM. The 3 features are audio waves, MFCCs, and Mel Spectrograms.

33 There are some of the interesting findings along the way. 1. Simplicity beats complexity. The simplest  
 34 model is a 3-layer MLP. It took the least time to train yet it has the best result of an accuracy of 89%  
 35 whereas LSTM took the longest time but has the worst result of 60%. 2. Hyperparameters are very  
 36 important. Some of them leads to a terrible result not better than random guessing. 3. Preprocessing  
 37 is so important. When doing MLP, non-scaled MFCCs led to a bad result below 15%. 4. Transfer  
 38 learning is so powerful that it saves a lot of time to train complex models. Pretrained models can be  
 39 used to extract features that can be used for the problem at hand. 5. VGG16 is for image classification  
 40 yet it can be used for the image representation of audio. In other words, the features extracted by an  
 41 image model is relevant to audio classification. However, it performs less than those audio dedicated  
 42 models, e.g. VGGish and YAMNet. After all, VGGish and YAMNet are trained against audios and  
 43 VGG16 is trained against images. 6. When got a bad result not better than random guessing, it could  
 44 be a data problem or a hyperparameter problem. 7. A CNN model's weight file is much bigger than  
 45 other, in this case, VGG16. It's because of the great number of weights, which is due to the fact that  
 46 each pixel is regarded as a feature in CNN. 8. When the accuracy is 100% and the val-accuracy is  
 47 70%, it means overfitting, so more anti-overfitting should be applied like dropout layer. Also, the  
 48 complexity of the model should be reduced, maybe by reducing the number of layers or units of a  
 49 layer. 9. If the train-accuracy is lower than val-accuracy, it feels like the model is not complex enough.  
 50 Maybe add more layers or units. 10. If all the units number goes to the maximum in a hypertuning  
 51 process, it might indicate an underfitting. 11. For VGG16, I first followed some tutorial online, but  
 52 got poor results. The reason is I cropped the images to 150 by 150 from the top left corner, during  
 53 which more than half of the information was lost since the original size is 432, 288. After applying  
 54 the full image, the result got much better. 12. tanh should be used with LSTM so that the cuDNN  
 55 kernels can be used to speed up the training process significantly. 13. Memory is so important to  
 56 deep learning training. Colab crashed a lot due to limit of memory. VGG16 needs the most memory  
 57 since it has a lot of learnable weights. 14. When I got an accuracy of 99%, it looked so suspicious  
 58 that I examined the codes several times and I found out that the reason was I had not initialed a new  
 59 model instance for it. Simply put, I used the fitted model that had seen the training and testing data so  
 60 even the validation accuracy is as high as a 99%.

61 The structure of the document is as follows. In section 2, related work around music genres  
 62 classification is introduced and compared with my approaches. In section 3, a detailed description  
 63 of the project is given with respect to the dataset used, preprocessing, baseline approach, hyper-  
 64 parameter tuning, and etc. In section 4, the results of these experiments are reported and analysed. In  
 65 section 5, an overview of the main findings is shown. Also, proposals regarding how the models can  
 66 be improved in the future are made. References lied in the section 6.

## 67 **2 Related work**

68 In 2002, Tzanetakis and Cook presented music genre classification as a pattern recognition task[1].  
 69 In that work, the authors assessed the classification using acoustic features extracted from the sound  
 70 on a dataset with 1000 music pieces labeled according to 10 musical genres. Instead, the MLP in this  
 71 document uses MFCCs to train a MLP. A baseline is also created in a similar way for this document.

72 Gwardys and Grzywczak [2] used a CNN trained on the Large Scale Visual Challenge (ILSVRC) as a  
 73 feature extractor. In their case, they did not have enough data to train a classifier, but enough data was  
 74 available in another domain of interest where the data could be in a different feature space or follow a  
 75 different data distribution. Similar to my work, I used VGG16, which is a visual model, to extract  
 76 features for processing audio. Pan and Yang [3] showed that this knowledge transfer, also known  
 77 as transfer learning, if done successfully, would greatly improve the performance of the learning  
 78 algorithm hence avoid expensive data-labeling efforts.

79 Sigitia and Dixon [4] trained a neural network with rectified linear units (ReLU) using Stochastic  
 80 Gradient Descent (SGD). They used the activations of the hidden layers of the neural networks as  
 81 features and trained a Random Forest classifier on top of these features to predict the classes. Similar  
 82 work in this document that use batch Stochastic Gradient Descent and applied ReLU on hidden layers

83 and trained MLP, CNN, and RNN on the top of the pretrained VGGish and YAMNet. Thus, transfered  
84 learning is also applied.

### 85 3 Experimental setup

86 The GTZAN dataset[5] is used for all models. It has 1000 30-second music pieces classified into 10  
87 classes. Each class has 100 audio pieces. The distributed dataset consists of three parts. The audio  
88 file in the wav. format. Mel Spectrograms in png. format. And the extracted MFCCs features in CSV,  
89 which have been summarized by taking the mean and std. All three forms of data are used. They are  
90 used as follows.

models	MLP	VGG16	LSTM	YAMNet	YAMNet&LSTM
data features	MFCCs	Mel Spectrogram	MFCCs	audio wave	audio wave

Figure 1: The implemented models and the data feature used

91 Each model is comprehensively hypertuned using KerasTuner[6] in terms of model architec-  
92 ture(whether a particular layer is included, e.g. Dense, LSTM, Dropout), layer units, learning  
93 rate, activation function, dropout rate, learning rate, batch size, and epochs. The best epochs is chosen  
94 separately using early stop callbacks, i.e. stopping the training before the result is not improving. Dif-  
95 ferent preprocessing methods are compared for MLP and VGG16. Hyperparameters are chosen based  
96 on a validation set, which is split from the training set. The weights from the best hyperparameter set  
97 is saved.

98 The preprocessing on the label data is one hot encoding since we are doing multiple classification  
99 problem here. The preprocessing on the training data varies in accordance with concrete models.

100 tanh activation function is used for LSTM and the hybrid of LSTM&YAMNet for optimised training.

#### 101 3.1 MLPs

102 The first model MLP, which contains 3 layers and takes summarized MFCCs(means and stds) as  
103 input features. The summarized MFCCs is included in the distributed dataset in two versions. One is  
104 for 30-second audio. The other is for 3-second audio. 3-second version is chosen since it led to better  
105 results. The MFCCs in the csv is not normalized. It was normalized before used.

106 MLP is chosen to be the baseline model because it is simple and effective.

#### 107 3.2 VGG16

108 The second model of a CNN model built and pretrained by Google, VGG16. Transfer learning is  
109 applied here because the CNNs are usually hard to train both in terms of computing power and dataset  
110 scale.

111 The data used is Mel Spectrogram, which is included in the distributed dataset in the png format.  
112 The images are converted into RGB, then three preprocessing methods are tried on them. 1) no  
113 preprocessing; 2) normalized by 255; 3) use a function preprocess\_input, which is in the package of  
114 tensorflow.keras.applications.vgg16.

115 On top of the pretrained VGG16, three or four(depend on a hyper parameter) Dense layers are applied.  
116 The first two layers have the sizes up to 3072(try to match the size of the VGG16 features), and then  
117 a layer with size of a range from 32 to 1024.

#### 118 3.3 YAMNet

119 Three Dense layers are seated on top of the pretrained YAMNet model.

120 The input is the audio wave files. 30-second audio pieces are loaded and then each is split into 10  
121 3-second smaller pieces to produce more data. At this stage, some pieces have different sizes with  
122 the majority. Therefore, a trimming process is applied here to include only the first 62 frames. Then,  
123 the data is summarized by take its mean and std across the time axis, after which the data is fed into  
124 the Dense layers.

### 125 3.4 LSTM

126 The input data is MFCCs. But the MFCCs csv. file in the given dataset can not be used because they  
127 are already summarized across time. So, the audio wave files are loaded, and split. After clearing  
128 some malformed samples, MFCCs features are generated from the 3-second audio pieces. At last,  
129 three LSTM layers are applied to these data.

### 130 3.5 Hybrid of LSTM&YAMNet

131 On top of YAMNet features, instead of putting MLPs, LSTMs are put on it. Specifically, Audio  
132 waves are used to generate YAMNet features. Then, three LSTM layers are applied on the extracted  
133 YAMNet features. Data preprocessing is the same with YAMNet.

## 134 4 Results

135 The evaluation phase involves splitting the data twice. The first time splits the data into training and  
136 testing datasets(The ratio is 4 : 1). The testing dataset will not be used until the very end to test the  
137 final generalizing ability of the model. The training data is then further split into two parts. One is for  
138 training and the other is for validating(The ratio is again 4 : 1). Validating is for finding the best set  
139 of hyper-parameters.

140 Early stop was used to determine the best epochs. With this technique, the training process stops  
141 itself when it sees no further improvement for the validating accuracy in a number of iterations. We  
142 need to set a patience value for this. The greater the value, the more epoches will be run.

143 The measure used here is accuracy, which describes the rate of correctly classified samples. In this  
144 case, it means how accurate the model is when classifying the samples. A value of 80% means it got  
145 80% of the sample correctly classified and the rest 20% incorrectly classified. An accuracy of 10% in  
146 our case means the model does not perform any better than random guessing.

147 For the models that are learned on 3-second audios features, the model can classify a genre by  
148 listening only to any 3-second piece of the audio. However, if that 3-second piece happen to be blank  
149 or silent, the classification would not be accurate. For the models that are learned on 30-second  
150 audios features, the model can classify a genre by listening to a 30-second piece of music.

151 executions\_per\_trial is set to 2 so that two executions are run for each trial instead of 1, so that the  
152 result is more stable and representative.

### 153 4.1 MLP

154 The MLP model reached the best result of 88.9% accuracy. The best set hyperparameters is as  
155 follows.

156 This model has the least amount of trainable weights and took the least amount of time to search the  
157 best hyperparameters and to train the model. Yet, it has the best result of a 88.9% accuracy. This  
158 model is regarded as the baseline model because of its simplicity.

159 I only did 5 epochs when doing hypertuning. As we can see above, both are increasing but the  
160 validation accuracy is flatter from 1 to 3 epochs. Two lines in each graph are the two executions for  
161 each trial. The two executions performed similarly in training but more differently for validating.

```

'activation': 'relu', 'units_1': 448, 'dropout_1': False,
'units_2': 672, 'dropout_2': False, 'units_3': 32,
'dropout_3': False, 'learning_rate': 0.001, 'batch_size': 12,
'epochs': 11

```

Figure 2: The best set hyperparameters for the MLP.

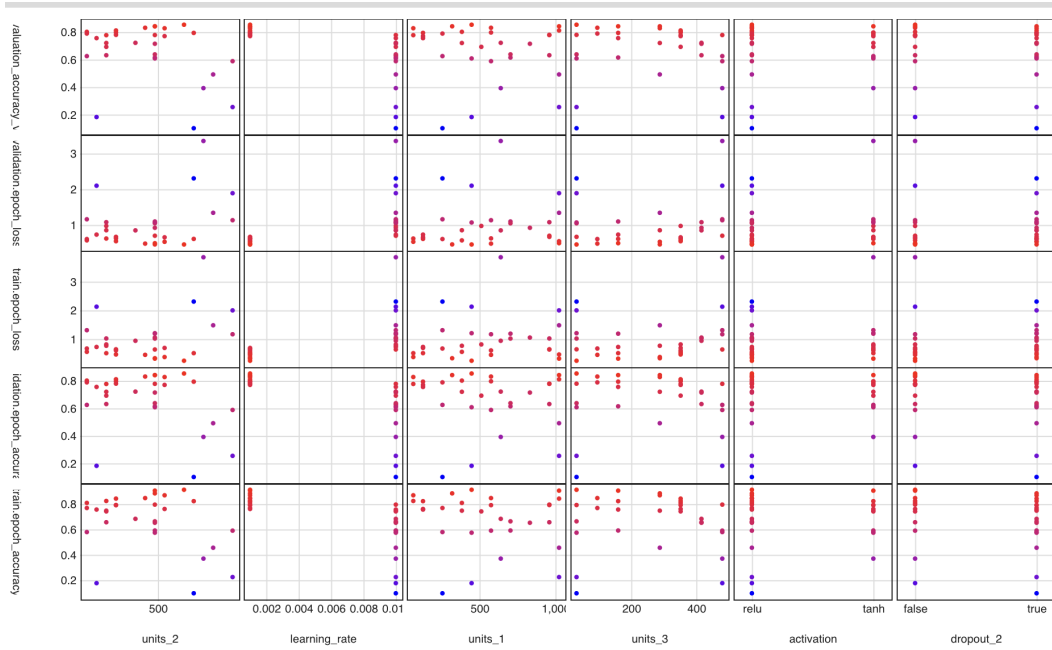


Figure 3: Part of the hyper parameter search space and result for the MLP.



Figure 4: Train and validation epoch accuracy

## 4.2 VGG16

The pretrained VGG16 is utilised to give a result of 65.5% accuracy on the following set of hyperparameters.

The result is much less than the baseline model. It took a considerable amount of memory to train due to the number of learnable weights. After VGG16 is an image model. The features it extracts is

```
'activation': 'relu', 'units_1': 3072, 'dropout_1': True,  
'units_2': 1536, 'dropout_2': False, 'units_3': 288,  
'dropout_3': False, 'has_forth_layer': False, 'dropout_4': False,  
'learning_rate': 0.001, 'units_4': 736, 'batch_size': 28,  
'epochs': 6
```

Figure 5: The best set hyperparameters for the VGG16.

167 for classify images not sounds. 65% is much better than random guessing, 10%. Therefore, VGG16  
168 features are helpful with audio classification in some way.

#### 169 4.3 YAMNet

170 YAMNet is a pretrained audio model, which is good at extracting audio features. This model reached  
171 a result of 84.9% of accuracy on the following hyperparameters.

```
'activation': 'relu', 'units_1': 1024, 'units_2': 512,  
'units_3': 128, 'dropout': True, 'learning_rate': 0.001,  
'batch_size': 52, 'epochs' : 5
```

Figure 6: The best set hyperparameters for the YAMNet.

172 This model is a bit less than the baseline model but close and much better than VGG16. After all, this  
173 extracts audio features instead of image features in the case of VGG16.

#### 174 4.4 LSTM

175 LSTM is a RNN, which is good at working with time series. This model reached a result of 77.3% of  
176 accuracy on the following hyperparameters.

```
'units_1': 416, 'units_2': 32, 'dropout_2': True,  
'units_3': 48, 'dropout_rate2': 0.2, 'learning_rate': 0.001,  
'batch_size': 12, 'epochs': 16
```

Figure 7: The best set hyperparameters for the LSTM.

177 This model is a less than the YAMNet model but still has a good accuracy of 77%, which is much  
178 more than VGG16. This model is really slow to train and very computing demanding.

#### 179 4.5 Hybrid of YAMNet&LSTM

180 The hybrid model has a result of 84.2%, which is significantly better than LSTM alone. The hybrid  
181 did show some advantage over its component. However, it achieved a similar result to the other  
182 component, YAMNet.

### 183 5 Conclusion & Future work

184 By exploring MLP, CNN, RNN, and transfer learning technique, music genres classification problem  
185 is solved by achieving high accuracies on the five models. MLP with MFCCs is the simplest yet

```
'units_1': 480, 'dropout_rate1': 0.15, 'units_2': 128,
'has_additional_layer': False, 'dropout_2': False,
'units_3': 80, 'dropout_rate2': 0.25, 'learning_rate': 0.001,
'units_additional': 64, 'batch_size': 44
```

Figure 8: The best set hyperparameters for the LSTM.

186 the best one. VGG16 is not for audio features but did okay. YAMNet is a fine model for audio  
 187 features and gave 85% of accuracy. LSTM worked well but took too much time to train. A hybrid of  
 188 LSTM&YAMNet worked better than pure LSTM but not better than pure YAMNet.

## 189 5.1 Findings

190 There are two additional findings that I would like to add here.

191 The first one is that there is really not that one architecture fit for all. The choice of architecture is  
 192 very open. It is not that more complex model is better. It is not that fancier models can beat the  
 193 simple ones. It is not true in the machine learning world.

194 The last one is how easy and cost-effective it is for an individual to build a decent deep learning  
 195 model to address a real world problem nowadays. That might be the output of the development of  
 196 science.

## 197 5.2 Future work

198 The champion among this models is the 3-layer MLP, which is the baseline model, yet, it discarded  
 199 most of the sequential information. LSTM that respects the sequential information should have better  
 200 performance yet the YAMNet and LSTM hybrid devised has achieved 85%, which is lower than 89%  
 201 achieved by MLP. One focus of future work might be experimenting and tuning the performance of  
 202 RNN models and its hybrid. Hopefully, better results can be achieved.

203 Try more anti-overfitting techniques, such as regularization, data augmentation, ensemble methods,  
 204 and cross-validation.

205 At last, there might be a defect with CNNs processing Mel Spectrograph because unlike photos,  
 206 vertical translation should not be implied on the Mel Spectrograph. The reason is different vertical  
 207 lines representing different pitch of sounds when the same pattern found across verticle lines. The  
 208 same shape in the higher pitch could represent totally unrelated things to that in the lower pitch.  
 209 Maybe, 1D CNN could help here because there is only one axis of translation.

## 210 References

- 211 [1][5] G. Tzanetakis, P. Cook Musical genre classification of audio signals IEEE Trans. Speech Audio Process.,  
 212 10 (5) (2002), pp. 293-302
- 213 [2] G. Gwardys, D. Grzywczak Deep image features in music information retrieval Int. J. Electron. Telecommun.,  
 214 60 (4) (2014), pp. 321-326
- 215 [3] S.J. Pan, Q. Yang A survey on transfer learning IEEE Trans. Knowl. Data Eng., 22 (10) (2010), pp.  
 216 1345-1359
- 217 [4] S. Sigtia, S. Dixon Improved music feature learning with deep neural networks IEEE International Conference  
 218 on Acoustic, Speech and Signal Processing (2014), pp. 6959-6963
- 219 [5] Karen Simonyan and Andrew Zisserman Very Deep Convolutional Networks for Large-Scale Image Recog-  
 220 nition karen,az@robots.ox.ac.uk

- 221 [5] Simonyan, K., & Zisserman, A. Very deep convolutional networks for large-scale image recognition. In  
222 Proceedings of the International Conference on Learning Representations (ICLR) (2015). Retrieved from  
223 <https://arxiv.org/abs/1409.1556>  
224 [6][https://keras.io/keras\\_tuner/](https://keras.io/keras_tuner/)