

Audio Classification using Convolutional Neural Networks

Aryan Singh

Email: aryansingh080501@gmail.com

Omer Al Sumeri

Email: omeralsumeri@gmail.com

Abstract—The objective of this project is to achieve precise audio classification across spoken digits from 0 to 9 utilizing a variety of Artificial Neural Network computation models. Specifically, Deep Convolutional Neural Networks (DCNN), MobileNet, ResNet-50, Inception V3, and YOLO (You Only Look Once) are employed for this classification task. Each model brings distinct capabilities and efficiencies to the table, contributing to a comprehensive exploration of audio classification techniques. Through rigorous experimentation and evaluation, this project aims to identify the strengths and weaknesses of these models in accurately recognizing and categorizing spoken digits.

Index Terms—Audio Classification, artificial neural networks, classification

I. INTRODUCTION

This project involves the use of various deep neural networks for audio classification as well as the preprocessing method of converting the .wav audio files into spectrograms. The different types of neural network models used were DCNN and then the results of this method were compared with the results of pretrained models like Mobilenet, ResNet-50, Inception V3 and YOLO. Python programming language was used for this project with imported deep learning libraries like TensorFlow, Keras and Scikit-learn and data manipulation libraries like Numpy and matplotlib.pyplot.

A. Preprocessing

In order to classify the audio files (.wav) using the various neural networks, they have to be converted into a form that could be spatially structured, similar to images. Therefore the audio files have to be converted to spectrograms, which capture both the time and frequency domains of an audio signal. The Spectrograms, as shown in Figure 1, depict a 2D representation of how the signal's frequency changes over time with the X axis representing time and Y axis representing the frequency at that instant. Using spectrograms makes it easier to perform dimensionality reduction, therefore allowing processes like DCNN to process and extract meaningful features.

B. DCNN

Deep Convolutional Neural Networks are a type of Artificial Neural Network designed for processing and analyzing data such as images and audio/video files. They consist of multiple layers, including convolutional layers, pooling layers, and fully connected layers. The convolutional layers detect spatial patterns and features in the input data, which are spectrograms in our case, and produce feature maps using multiple filters.

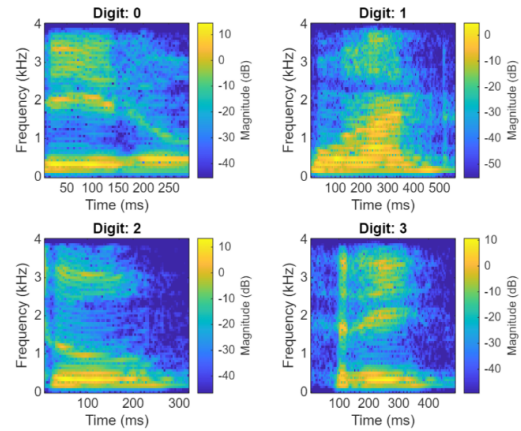


Fig. 1. Spectrograms for different digits (0, 1, 2, 3)

This process of extracting important features from the input data using filters is known as convolution. These filters/kernels slide over the input data and extract the local features. Each filter detects a specific pattern or feature such as edges, textures or shapes at different spatial locations of the input data (different frequencies at different instances in the case of spectrograms). The second layer of DCNN are the pooling layers, which follow the convolutional layers. These layers reduce the spatial dimensions of feature maps while retaining the most important features. The Max-Pooling layers select the maximum values within each pooling window, therefore preserving the most dominant features of the feature maps. This layer is essential to reduce computational complexity, memory requirements and prevent overfitting by reducing the spatial dimensions of the data. The final layer of the DCNN architecture is the fully connected layer or the dense layer. These layers perform data classification and regression tasks on the features extracted from the previous 2 layers. These layers connect each neuron from the previous layer to every neuron in the current layer to form a fully connected structure. These layers can perform decision making based on the features extracted from the previous layers which allows the network to make predictions or classifications. Functions like softmax (used by our team) are used by these layers to perform classification using probabilistic predictions i.e. converting raw scores into probabilities, therefore allowing the neural networks to make predictions for multi-class classifications, based on these probabilities. Figure 2 shows the entire DCNN

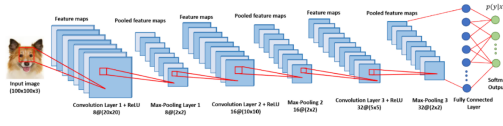


Fig. 2. Deep Convolutional Neural Network Architecture

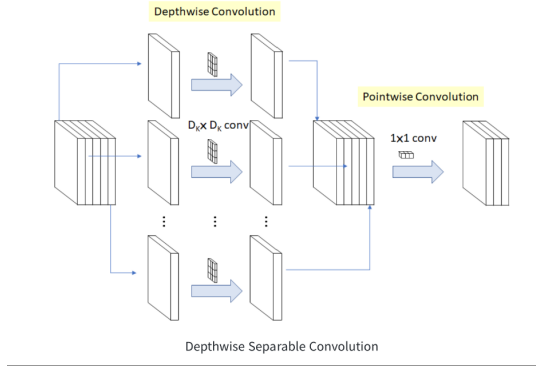


Fig. 3. Depthwise Separable Convolution in MobileNet

architecture, representing how these 3 layers work together and perform convolution and pooling, followed by classification.

C. MobileNet

MobileNet is a type of Convolutional Neural Network (CNN) architecture that is designed to be lightweight and efficient for image classification on mobile or embedded devices. This type of architecture is designed for devices with limited computational capacity since it minimizes the computational resources required for classification while ensuring high accuracy. MobileNet employs depthwise separable convolution, which splits the standard convolution into 2 separate layers: depthwise convolution and pointwise convolution. Depthwise convolution is a technique that applies a separate filter for each input channel, which means that there is a separate filter associated with each input for capturing spatial information specific to each channel. Pointwise convolution takes the feature maps created by the depthwise convolution filters and uses a 1X1 filter to combine all of the outputs, thereby reducing computational complexity while preserving the most important features. Figure 4 represents how Depthwise separable convolution is performed in MobileNet.

D. ResNet-50

ResNet-50 is a variant of the ResNet (Residual Network) architecture, a deep convolutional neural network. It is a powerful image classification model which uses residual connections to learn deeper architectures more effectively. Residual connections allow the network to learn a set of residual mappings and bypass layers for faster and more effective propagation. The ResNet architecture is divided into 4 main sections: the convolutional layer, the identity and convolutional block, and the fully connected layers. The convolutional layers are responsible for feature extraction from the input data. The

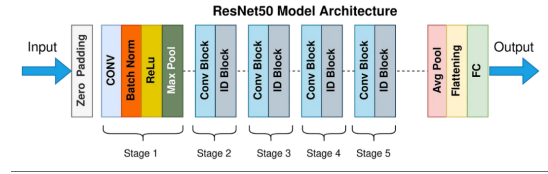


Fig. 4. ResNet-50 Architecture

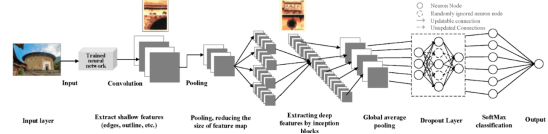


Fig. 5. Inception V3 Architecture

feature maps are parsed on by the convolutional layers to the max pooling layers, which reduce the spatial dimensions of the feature maps while preserving the most important features. The identity and convolutional blocks are responsible for processing these features and then the fully connected layers perform the final classification. The output of these layers goes into the softmax activation function which performs classification using probabilistic predictions for multi-class classification. Figure 4 represents the architecture of the ResNet-50 model.

E. Inception V3

Inception V3 is a state-of-the-art convolutional neural network architecture designed for advanced image classification. This architecture provides the perfect balance between accuracy and computational efficiency, which is why it is a very popular choice for computer vision applications. This model uses inception modules to capture multi-scale features efficiently by performing parallel convolutions with different filter sizes. The use of multiple filters and parallel convolution at each layer allows these modules to capture features at various levels of abstraction. This means that some filters are used to capture the finer details while others are used to capture the broader information, while working in parallel with each other. The feature maps processed by each parallel convolution layer are concatenated with each other and the final output has feature maps captured using different scales and filters, therefore representing information rich data. This model also involves various factorization techniques like performing 1X1 convolutions for dimensionality reduction before each parallel convolution path. Through this the model is able to reduce complexity and allow more efficient processing by the inception modules. Along with these layers, the architecture also has pooling layers for downsizing the spatial dimensions of feature maps. Figure 5 represents the architecture of Inception V3.

F. YOLO

YOLO (You Only Look Once) is a state-of-the-art object detection algorithm known for its incredible speed and ac-

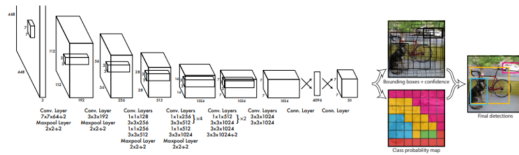


Fig. 6. YOLO Architecture

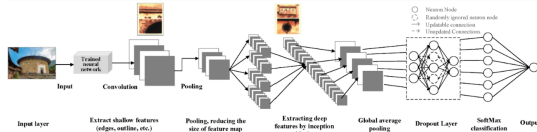


Fig. 7. ResNet-50 Architecture

curacy with which it classifies real time objects in images. The reason why YOLO is much quicker than other object detection algorithms is because it performs object detection in a single pass through an image rather than the multi stage approach that traditional methods take. In a single pass through a neural network, this model predicts the bounding boxes and class probabilities from the entire image. The learning approach that YOLO takes is that it divides the input image into a grid of cells and then predicts the bounding boxes and class probabilities for the objects that lie inside each cell. This network outputs a probability distribution over all possible classes for each bounding box. After all predictions are made for all cells, this architecture uses a post processing method called non-max suppression (NMS) to remove redundant or overlapping bounding boxes. This is done by NMS by selecting the most confident bounding boxes based on the predictions made and the rest of the overlapping boxes are removed to avoid redundancy. Figure 6 shows the architecture of YOLO, showing all the convolutional layers followed by the bounding boxes and class probability maps. Combining these 2 give the final detection results

II. PROCEDURE

A. DCNN

The development of the DCNN architecture was done with the use of various python libraries like “os” for file operations, “librosa” for audio processing, “numpy” for numerical computations, “tensorflow” and “keras” for machine learning and “matplotlib.pyplot” for data visualization. After all the libraries were specified, the preprocessing task of converting the audio files to spectrograms was performed. For this task all the parameters for the spectrograms and their labels were declared. Each audio file was read and its corresponding spectrogram was created using Short Term Fourier Transform (STFT), followed by truncation or padding to ensure a fixed length for each spectrogram. Once the input data was converted into spectrograms, the DCNN architecture was defined and the convolutional layers (3 in total) along with ReLU (Rectified Linear Unit) activation unit, and the max pooling layers were introduced. Finally the output of the max pooling layer was

parsed onto the dense layers with softmax activation function for accurate classification. Therefore the model was designed and the Adam optimizer was used for compiling. The data was then split into training and validation sets, where 80% of data was reserved for training and 20% for validation, followed by training being done for 32 batch size and 20 epochs. Once initial training was done, the code was optimized and a smaller learning rate of 0.0001 was defined to retrain the data. Once the training was done the accuracy vs epoch graphs are plotted, showing the model's performance on the validation set. The confusion matrix was also generated to further visualize DCNN's classification on the validation set.

B. MobileNet

The development of the MobileNet architecture also involves the use of all the same libraries as the ones used for DCNN architecture development. Once all the libraries were loaded in, the preprocessing of the raw audio files was done by converting them to spectrograms. This was done by using Short Term Fourier Transform (STFT), followed by truncation or padding to ensure a fixed length for each spectrogram, also similar to DCNN design. Once the preprocessing was completed, it is then splits the preprocessed data (spectrograms) into training and validation sets (80% for training and 20% for validation) The MobileNet architecture was then loaded using the Keras library and the different convolutional, pooling and dense layers were added. The first dense layer had 128 units with ReLU activation and the final one had 10 units with softmax activation for multi class classification (all 10 digits). The model was then compiled with the Adam optimizer and the training was done for 20 epochs with the validation being performed based on the validation set of the model. The results of training and validation were then represented on a graph, showing the accuracy of the model with increasing epochs. Finally the confusion matrix was also generated to represent the performance of the model for each class (each digit in our case). Different classes were shown, with true positives (along the diagonal), false positives, true negatives and false negatives.

C. ResNet-50

For the development of the ResNet-50 architecture, all similar libraries were used for loading audio files, processing them and then performing machine learning operations. The spectrograms were then generated using short term fourier transform and padding was done to ensure fixed length for each spectrogram. Once the preprocessing was performed, the spectrograms were then split into training and validations sets (80% for training and 20% for validation). The architecture for ResNet-50 was loaded from the Keras library to include the convolutional, pooling and dense layers. Then the model was compiled using the Adam optimizer and trained on the training set for 20 epochs while performing validation on the validation set. The results were then plotted on a graph to show the accuracy of the model with each increasing epoch. Finally, the confusion matrix was generated to represent the

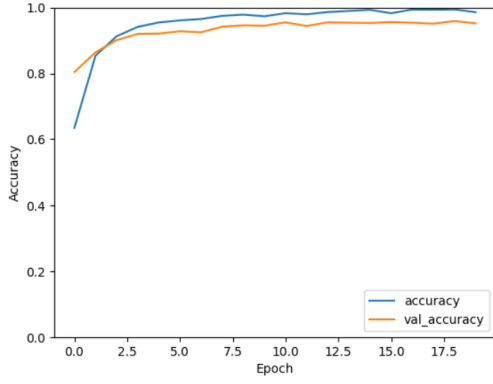


Fig. 8. Accuracy vs Epoch graph for DCNN using Adam optimizer

accuracy of the model in classifying each class (digit) in the validation set.

D. Inception V3

For developing the Inception V3 architecture, all the same libraries for loading audio files, array manipulation, deep learning and preprocessing of data were included. Similar to previous architectures, the audio files were loaded and converted to spectrograms using short term fourier transform, followed by padding or truncating them to ensure fixed length for each spectrogram. The data was then split into training and validations sets, with 80% data being reserved for training and 20% for validation. The Keras library was then used to load in the pretrained Inception V3 model, along with the convolutional, pooling and dense layers. After the model was loaded, it was compiled using the Adam optimizer and trained on the training set for 20 epochs while performing validation on the validation set. The results were then plotted on a graph to show the accuracy of the model with each increasing epoch. The confusion matrix was then generated to represent the accuracy of the model in classifying each digit in the validation set.

III. ANALYSIS AND DISCUSSION OF RESULTS

A. DCNN

For the custom DCNN model, two experiments were performed. One with the Adam optimizer and the other with a custom learning rate. Both models will be trained with a batch size of 32 and 20 epochs. Figure 7 contains an Accuracy vs Epoch graph for the custom DCNN model created with the Adam optimizer. After 3 epochs, the training and validation accuracy are around the 90% mark. After 6 epochs, the training and validation accuracy stabilize near the 95% range. The final validation accuracy is 95% for the custom DCNN model using the Adam optimizer.

In Figure 8, the Accuracy vs Epoch graph for the custom DCNN model with a 0.0001 learning rate. The smaller learning rate tells the model to be more thorough with the dataset in each epoch. The training accuracy and validation accuracy is above 90

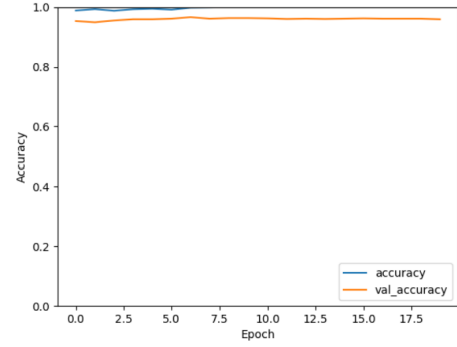


Fig. 9. Accuracy vs Epoch graph for DCNN using smaller learning rate

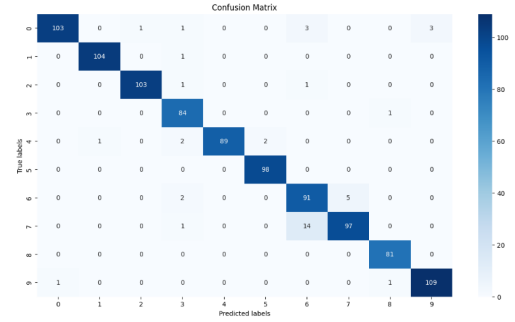


Fig. 10. Confusion matrix of DCNN model

Figure 9 contains the confusion matrix of the DCNN model with the smaller learning rate. The diagonal line of the matrix shows which classifications were predicted correctly. All classifications had at least above 80

B. MobileNet

For the low end model, MobileNet will be the model. The pretrained MobileNet model from tensorflow with the Adam optimizer will be used. One experiment will be conducted training the model using Mobile Net. The model will be trained with a batch size of 32 and 20 epochs.

Figure 10 is a plot of the Accuracy vs Epoch graph. After 5 epochs, the training and validation accuracy is above 85%. After 8 epochs, the accuracies stabilize out around the 85% mark. The final accuracy rate for the MobileNet model is 93%.

Figure 11 contains the confusion matrix of the MobileNet model. The diagonal line of the matrix shows which classifications were predicted correctly. All classifications had at least above 75

C. ResNet-50

For the midtier model, Resnet50 will be the model. The pretrained ResNet-50 model from tensorflow with the Adam optimizer will be used. One experiment will be conducted training the model using Mobile Net. The model will be trained with a batch size of 32 and 20 epochs.

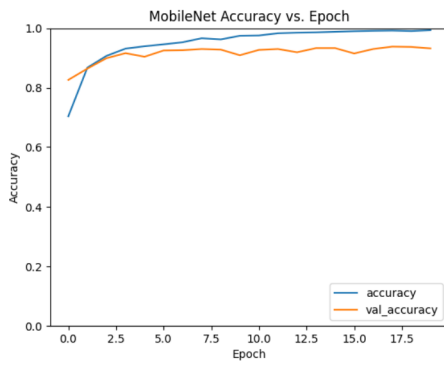


Fig. 11. Accuracy vs Epoch graph for MobileNet

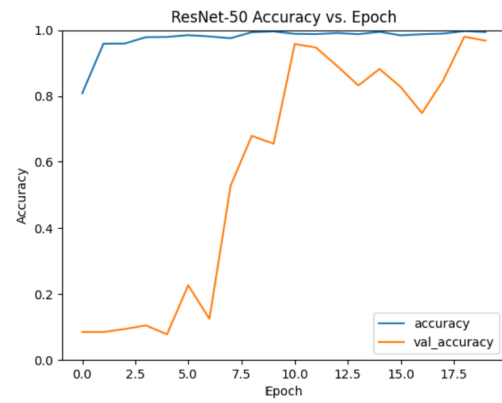


Fig. 13. Accuracy vs Epoch graph for ResNet-50

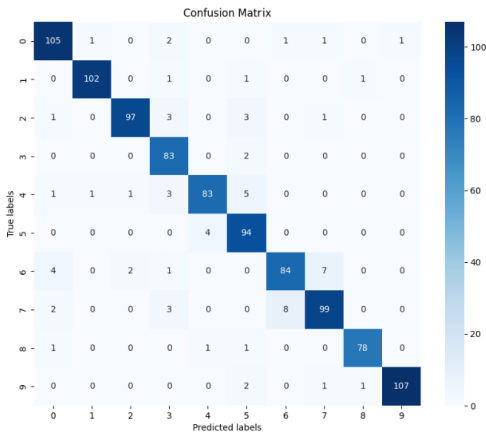


Fig. 12. Confusion matrix of MobileNet model

Figure 12 contains the Accuracy vs Epoch plot for the ResNet-50 training model. After 3 epochs, the training accuracy is high, above 90%. The validation accuracy is consistently below 20% until 7 epochs. From there, the validation accuracy is volatile. The reason for this could be overfitting, where the model tries to train the data too closely. This makes sense as Resnet50 is a mid tier model. At the end, the final validation accuracy is 93%.

Figure 13 contains the confusion matrix of the ResNet-50 model. The diagonal line of the matrix shows which classifications were predicted correctly. All classifications had at least above 75

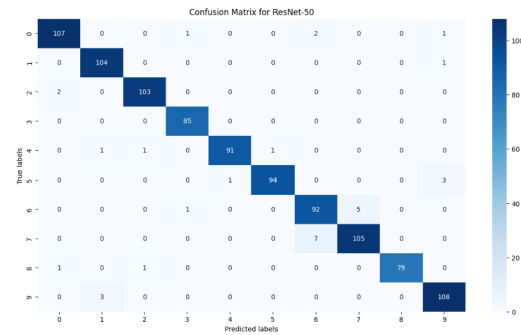


Fig. 14. Confusion matrix of ResNet-50 model

is just above around 70%. For both the training and validation accuracy were trending upward. The final validation accuracy is 90%.

Figure 15 contains the confusion matrix of the InceptionV3 model. The diagonal line of the matrix shows which classifications were predicted correctly. All classifications had at least above 80% accuracy. Most classifications hover around the 80% range. The worst performing class was the spoken digit 8, with an accuracy rate 81%. One interesting observation about

D. Inception V3

For the high end model, InceptionV3 is the model. The pretrained InceptionV3 model from tensorflow with the Adam optimizer will be used. One experiment will be conducted training the model using Mobile Net. The model will be trained with a batch size of 32 and 20 epochs.

Figure 14 contains the Accuracy vs Epoch plot for the InceptionV3 training model. Since Inception V3 is a high end model, the validation accuracy is volatile. Up until 5 epochs, the validation accuracy is around 60%. The training accuracy

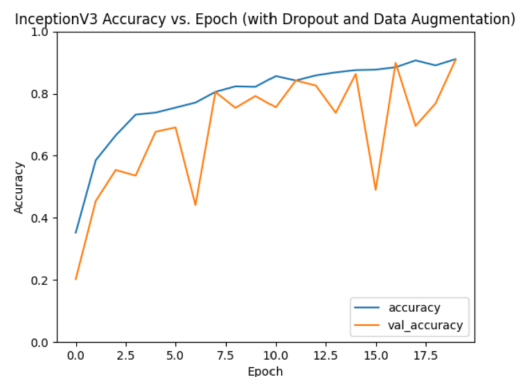


Fig. 15. Accuracy vs Epoch graph for InceptionV3

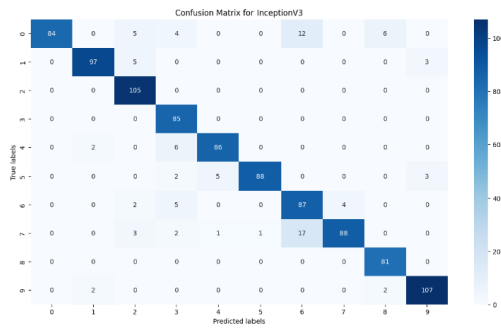


Fig. 16. Confusion matrix of InceptionV3 model

this confusion matrix is that spoken digit 7 was misclassified as spoken digit 6, 17 times.

E. YOLO

Results for YOLO were not achievable with the current hardware that was used to complete the training for the other models. Due to classifying .wav files, the .wav files need to be converted into spectrograms; the spectrograms would then be converted into a picture format. Even if the .wav file might be the same size, when they are converted into images, they are not the same dimensions.

For future work, the team would like to convert all .wav files into spectrograms, and the spectrograms will be converted to an image file. Once all the images are gathered, the largest image would be found and that would be the set size. All other images would be filled with black boxes. The team does not know how the accuracy will be affected; images of the same class will have different amounts of black boxes.

F. Analyzing all classification models

Table 1 contains the Validation accuracy and time to train for each classification model. The best performing classification model is the DCNN with a validation accuracy of 95%. The worst performing classification is InceptionV3 with a validation accuracy of 90%.

The classification model that took the shortest to train is DCNN, only taking 10 minutes. The classification model that took the longest is InceptionV3, taking 3.45 hours to train. The time to train for all the algorithms falls in line with the amount of layers each has. DCNN only had 3 layers. Each algorithm has more layers then the last so the time to train is longer.

Classification Model	Validation accuracy	Time to train
DCNN	95%	10 minutes
MobileNet 1	93%	10.66 minutes
ResNet-50	93%	2.22 hours
InceptionV3	90%	3.45 hours

TABLE I

VALIDATION ACCURACY AND TIME TO TRAIN OF CLASSIFICATION MODELS

CONCLUSION

In conclusion, this project was able to provide us with great insight into how spoken digits in audio files are classified using different deep neural networks like DCNN, MobileNet, ResNet-50, Inception V3 and YOLO. Along with learning about these models and their implementation, we were also able to understand which model works best for our dataset and gives the most accuracy and which model does not work well with spectrograms.

LINK TO DATASET USED

<https://www.kaggle.com/datasets/sripaadsrinivasan/audio-mnist?resource=download>