

Contents

1. Introduction 2

2. Methodology2

3. Experiments and Results 3

4. Code:13

Deep learning system for automatic speech recognition

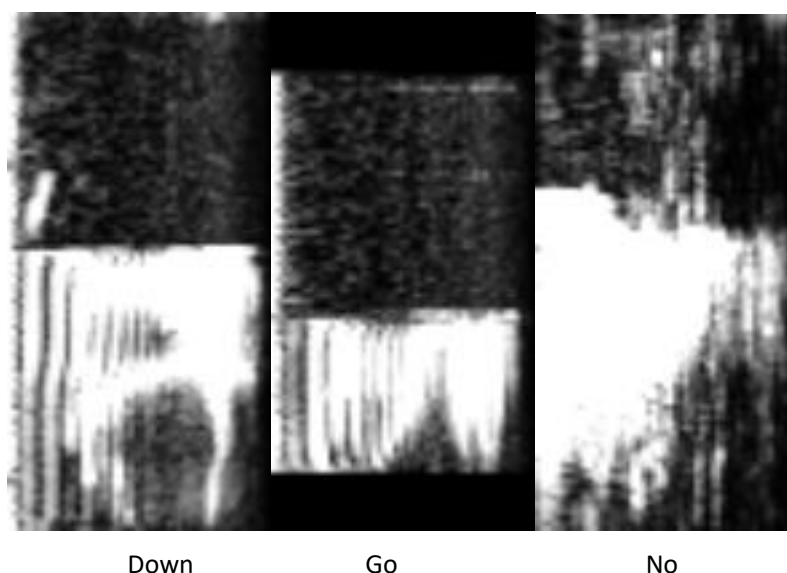
1. Introduction

Speech recognition is a fundamental task in natural language processing (NLP) and has become increasingly important with the rise of smart assistants and voice-controlled devices. Deep learning has shown great promise in improving the accuracy of speech recognition systems. This technical report presents an analysis of a deep learning-based approach to speech recognition, focusing on the design, implementation and evaluation of different neural network architectures and training strategies.

2. Methodology

The study used a large dataset of speech recordings, including both clean and noisy speech. The dataset was divided into training, and validation sets.

The input data has been pre-processed into spectrograms of speech 'images' like this:



speech waveforms in the time-domain and corresponding speech image spectrogram that forms the input data. There are 12 word classes, including unknown class and background class, the distribution shown in the plots:



Different neural network architectures, including convolutional neural networks (CNNs), recurrent neural networks (RNNs), and transformer networks, were trained using different training strategies, including supervised learning and transfer learning. For this study, I designed CNNs for speech recognition using Matlab for the following reasons:

- (1) Local feature extraction: CNNs are able to extract local features from speech signals, which can capture the patterns and characteristics of speech sounds. This is achieved by convolving a set of filters over small segments of the input signal to detect specific features, such as phonemes or syllables.
- (2) Translation invariance: CNNs are invariant to small shifts in the input signal, which can be useful in speech recognition where the same sound can be pronounced in slightly different ways. This property allows CNNs to capture the same feature regardless of its position in the input.
- (3) Parameter sharing: CNNs share parameters across the entire network, reducing the number of trainable parameters required for the model. This makes CNNs more efficient to train and less prone to overfitting, which is especially important in speech recognition where large amounts of data are required.
- (4) Hierarchical representation: CNNs are capable of learning hierarchical representations of speech signals, which can capture both low-level features (such as phonemes and syllables) and high-level features (such as words and phrases). This allows CNNs to model the complex relationships between different speech components and make accurate predictions.

The tasks accomplished in the technical report are:

- Design, implement and evaluate a basic deep learning classifier system to recognise spoken words from spectrogram image inputs and achieve an accuracy of $\geq 60\%$ on the validation data.
- Use a 2x2 grid search to perform hyperparameter tuning in the model of:
 - the number of convolutional layers.
 - the number of filters per layer.
- Reduce the computational complexity of the model using depthwise separable convolutions and 1x1 convolutions, and improve model generalisation.

3. Experiments and Results

3.1 Basic deep learning classifier system for speech recognition:

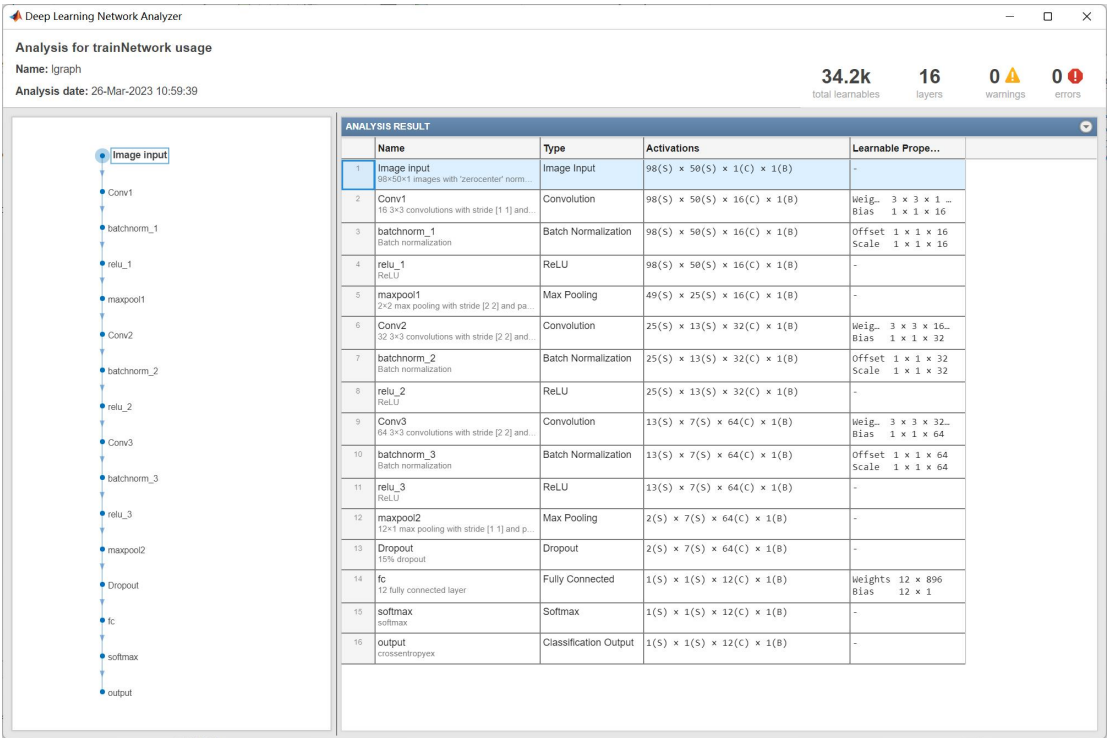
The model has 16 layers and 34.2k learnables:

- The input layer takes in an speech image with a size of 98 x 50 x 1.
- The first convolutional layer applies 16 filters to the input. each with a size of 3 x 3 pixels. The ReLU activation function is applied to the output of each filter, and batch normalize the data before next layer.
- The first max pooling layer reduces the size of the output from the first convolutional layer by applying a 2 x 2 window with stride 2 (i.e. the window moves 2 pixels at a time) and selecting the maximum value in each window.
- The second convolutional layer applies 32 filters to the output of the first max pooling layer, each with a size of 3 x 3 pixels and stride 2. The ReLU activation function is applied to the output of each filter, and batch normalize the data before next layer.
- The third convolutional layer applies 64 filters to the output of the second convolutional layer, each with a size of 3 x 3 pixels and stride 2. The ReLU activation function is applied to the output

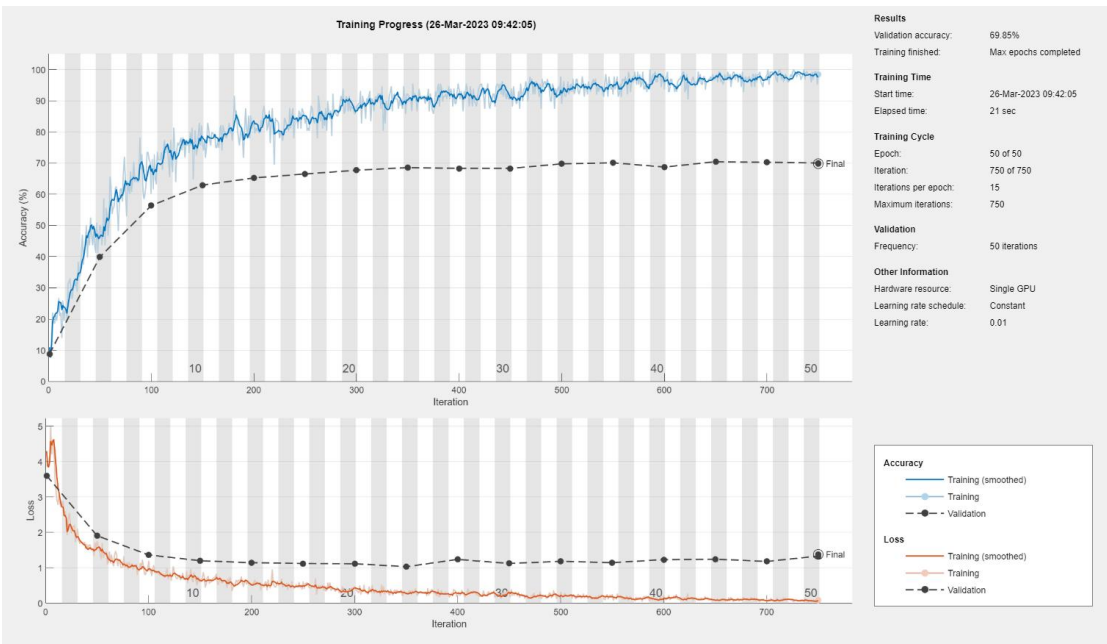
of each filter, and batch normalize the data before next layer.

- The second max pooling layer reduces the size of the output from the third convolutional layer in the same way as the first max pooling layer.
- The dropout layer makes some neural function stop working to avoid overfitting.
- The fully connected layer has 12 units.
- The output layer has 12 units (assuming a 12-class classification problem) and uses a softmax activation function to produce a probability distribution over the classes.

Network structure:



Training progress:

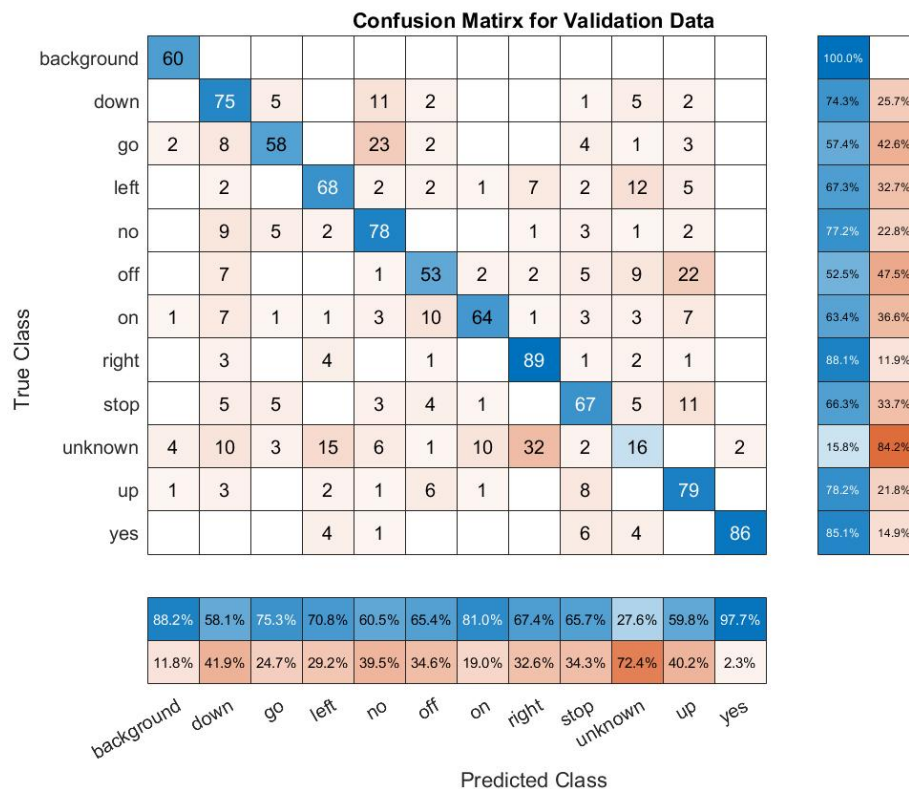


I trained this model on a single GPU like Nvidia RTX3060, some Hyperparameters are:

- Initial learning rate: 0.01
- Max epochs: 50
- Validation Frequency: 50
- Learning rate drop period: 20
- Learning rate drop factor: 0.1

The model achieved an accuracy of 69.85% on the validation data.

Confusion plot:



3.2 Grid search to perform hyperparameter tuning in the model of:

- the number of layers.
- the number of filters per layer.

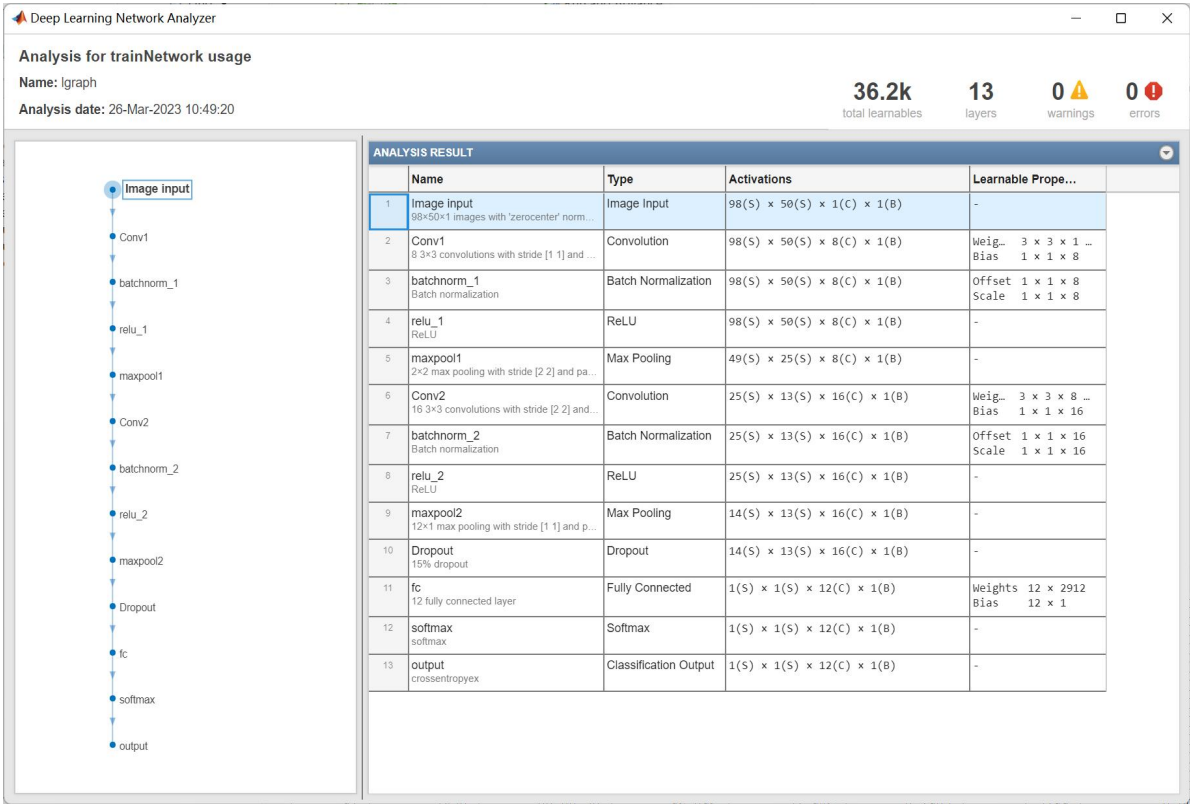
To avoid run into computational difficulties, I took 2x2 grid search as follows:

Num of conv layers: 2 Num of filters : 8, 16,.....	Num of conv layers: 2 Num of filters : 16, 32,.....
Num of conv layers: 3 Num of filters: 8, 16,.....	Num of conv layers: 3 Num of filters : 16, 32,.....

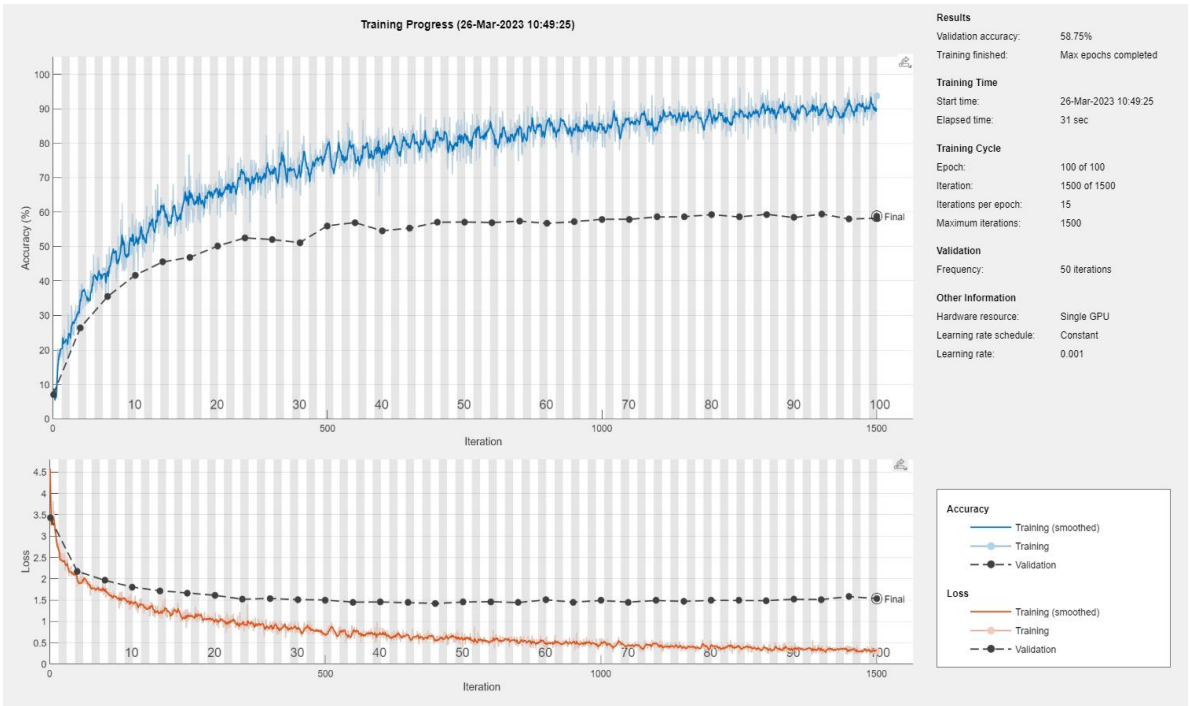
The number of convolution layers: 2

The number of filters:8, 16,.....

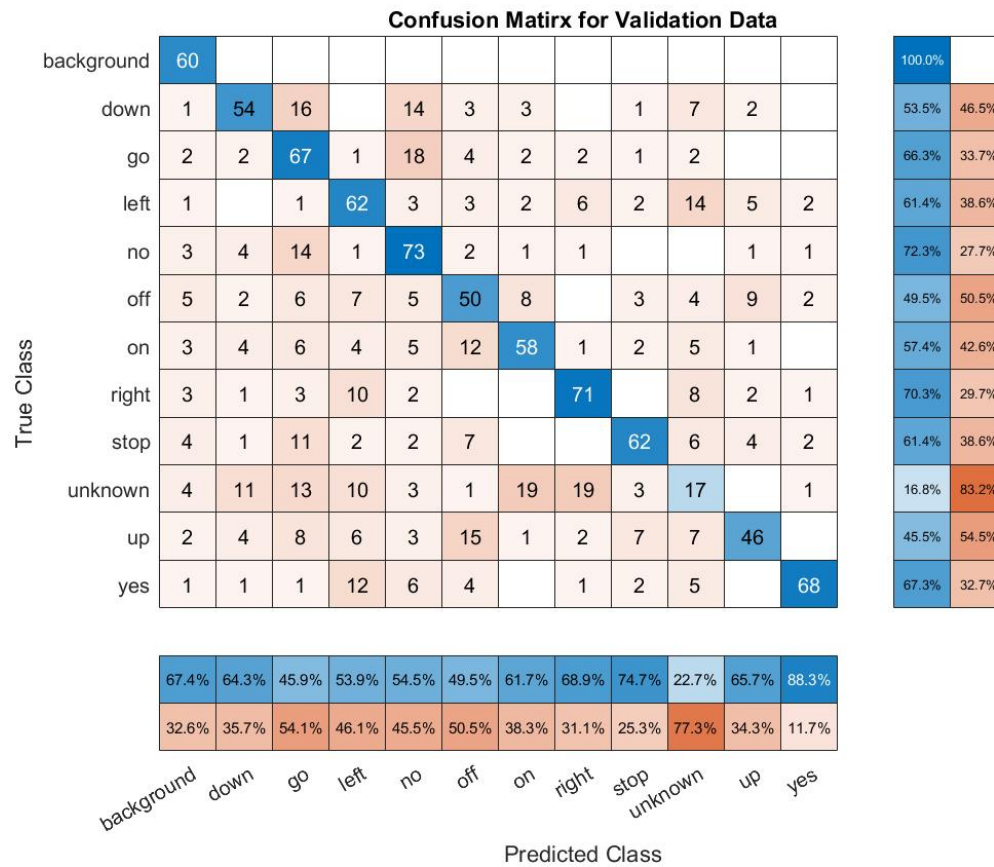
Network structure:



Training progress:



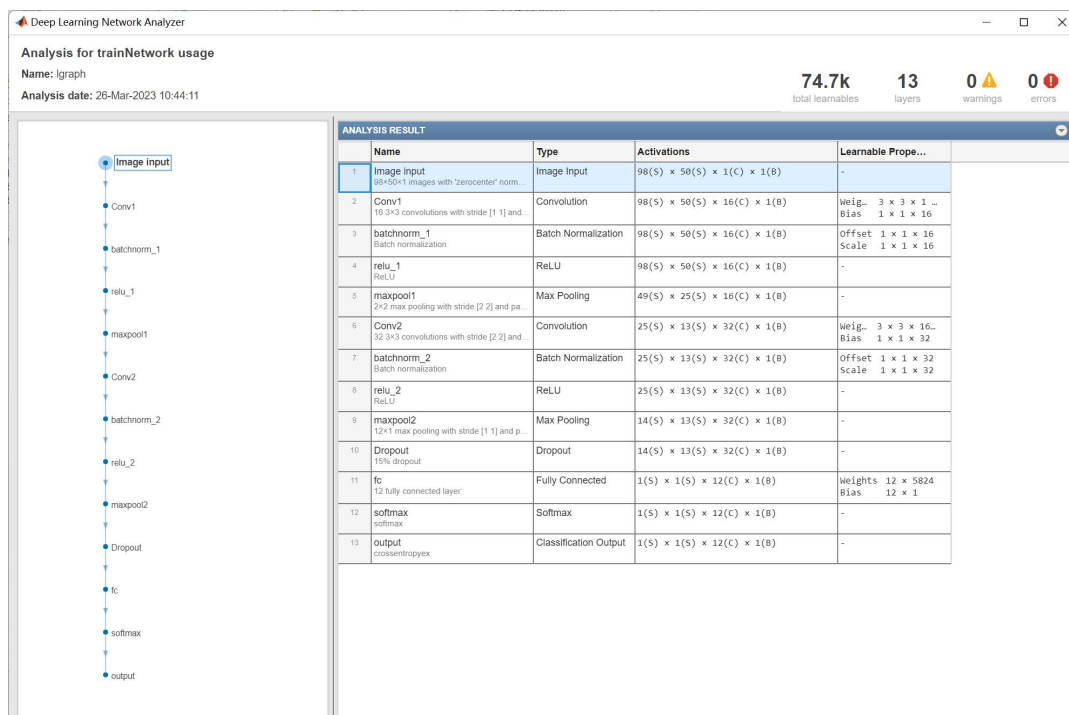
Confusion plot:



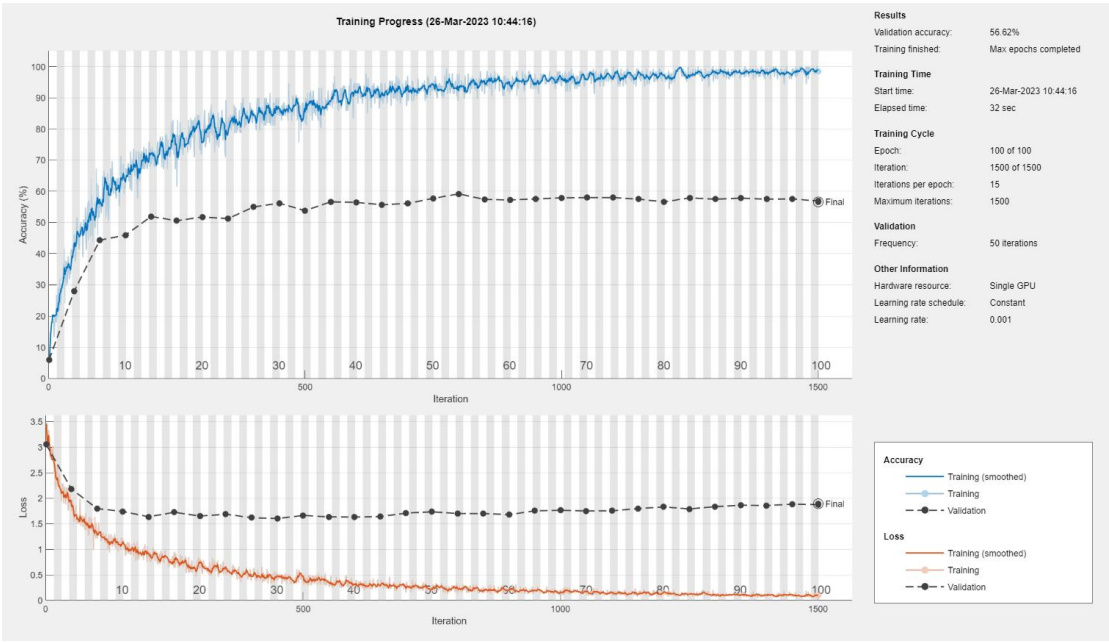
The number of convolution layers: 2

The number of filters: 16, 32, ...

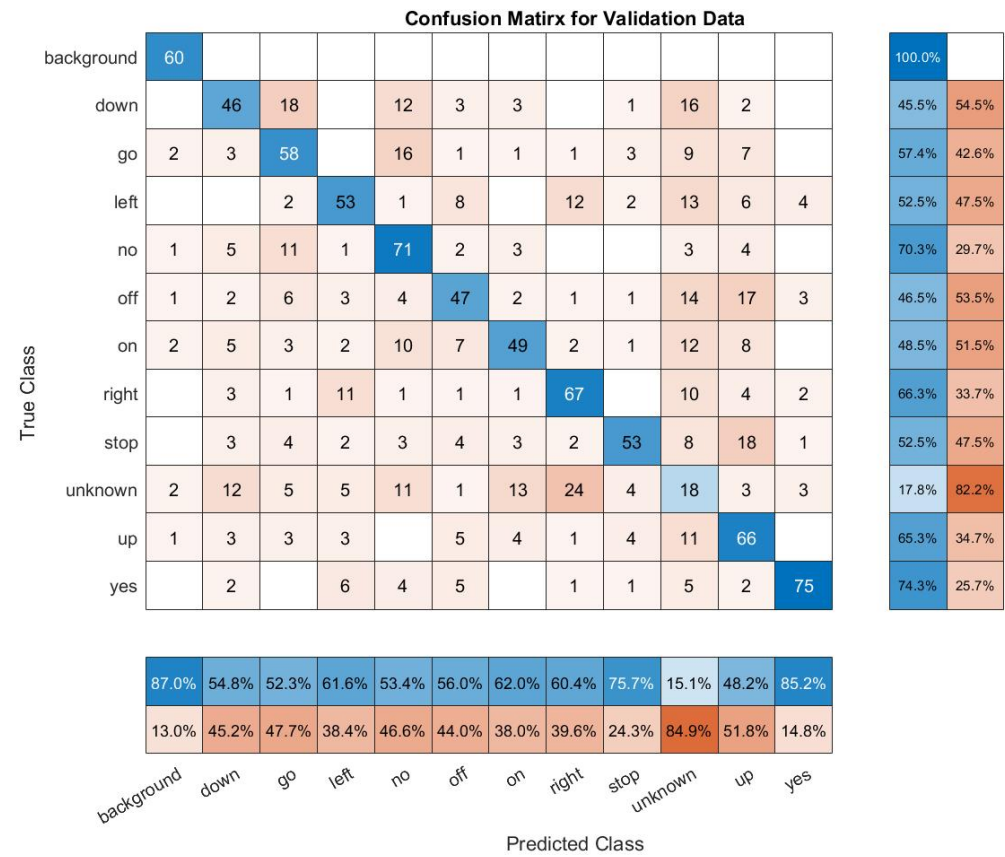
Network structure:



Training progress:



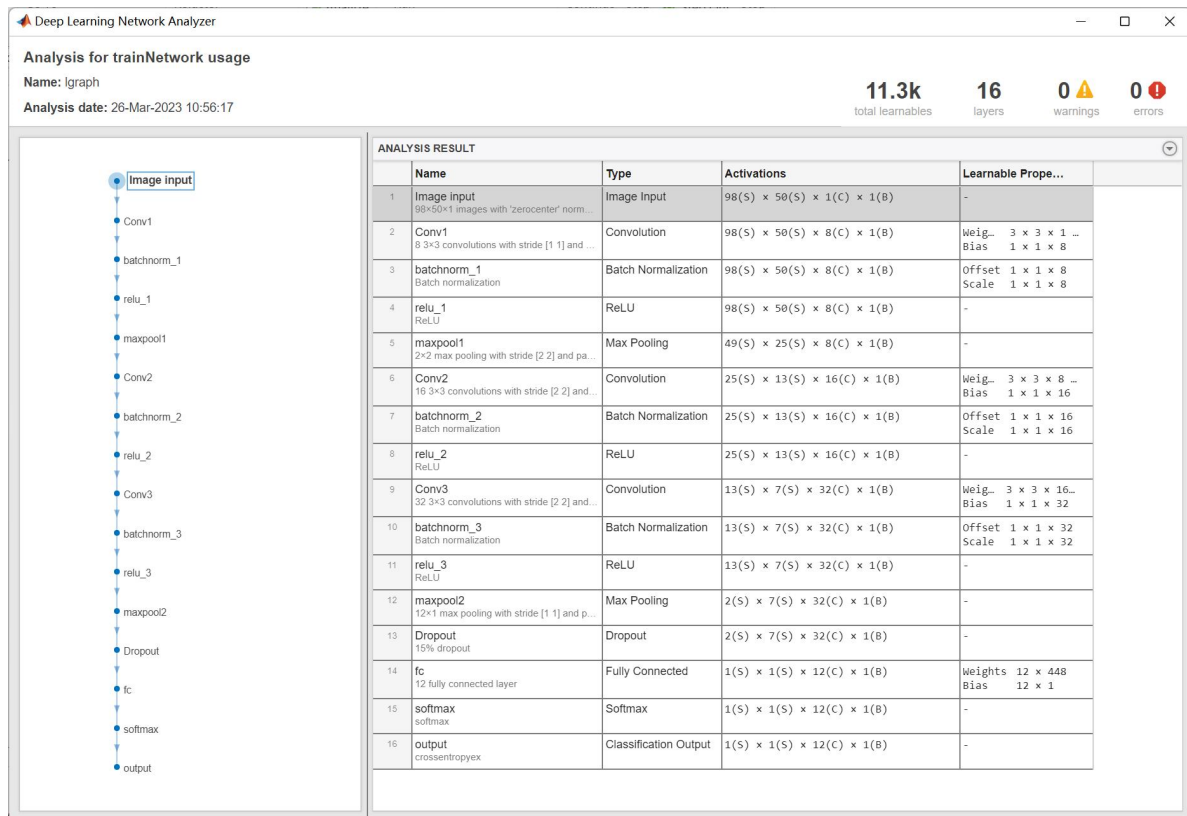
Confusion plot:



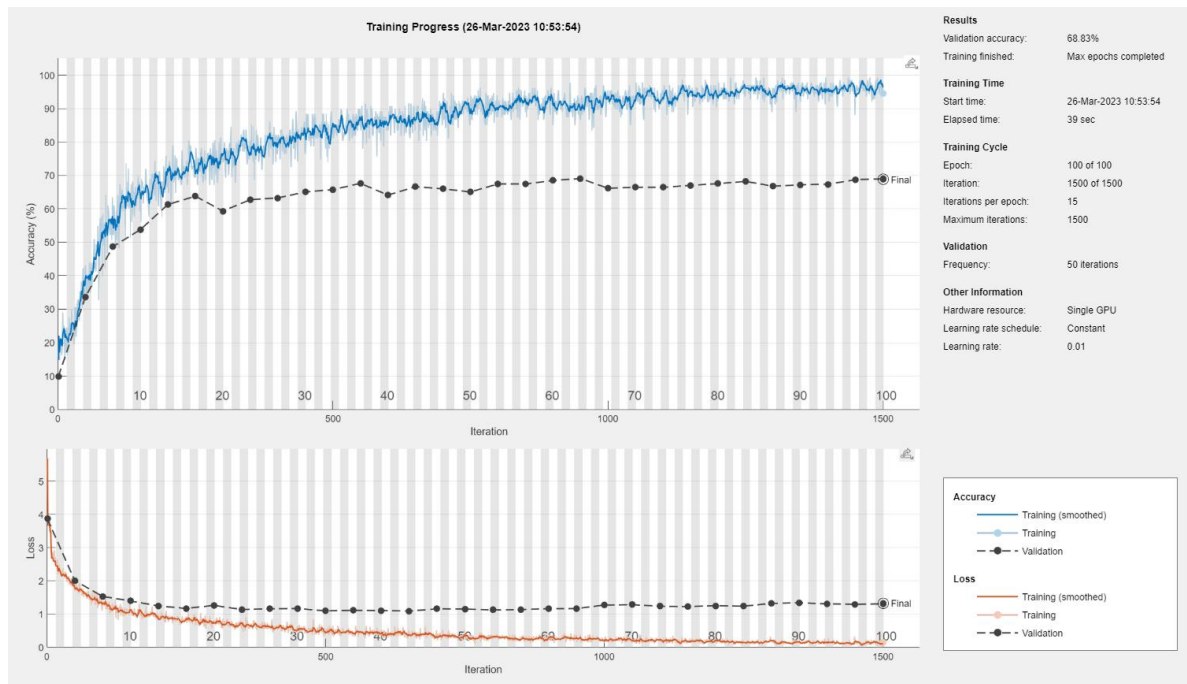
The number of convolution layers: 3

The number of filters: 8, 16,

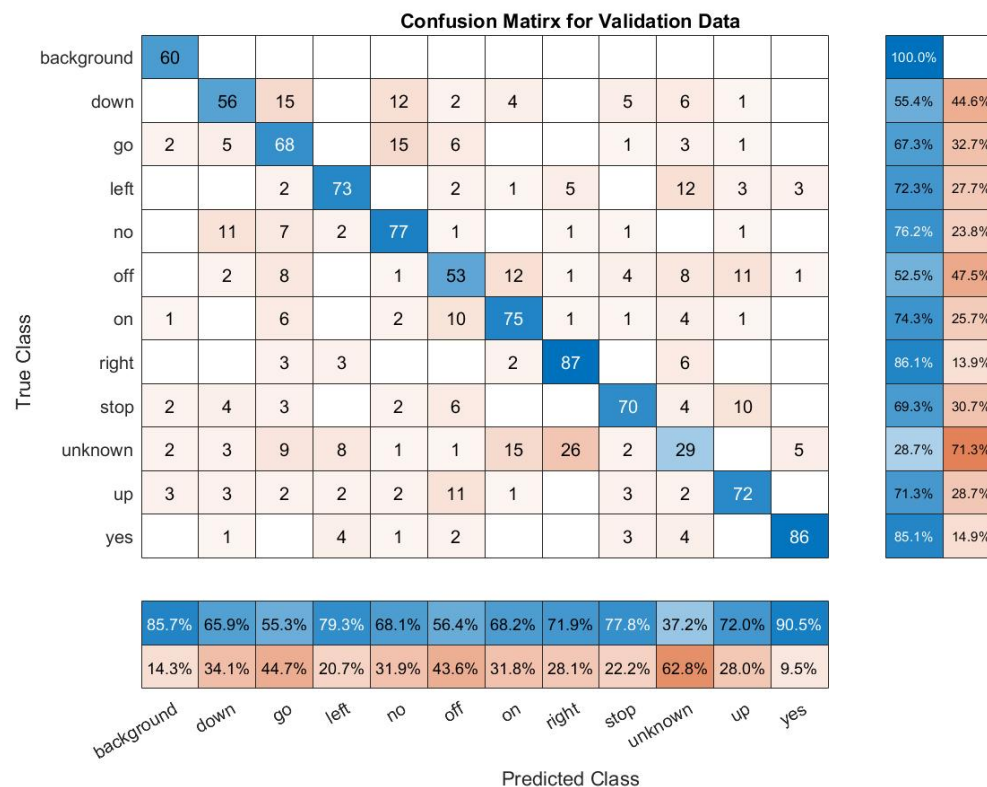
Network structure:



Training progress:



Confusion plot:



The number of convolution layers: 3

The number of filters of first convolution layer:16:

The details of this model are mentioned above(3.1).

3.3 Reduce the computational complexity:

In this study, I applied depthwise separable convolutions and 1x1 convolutions to reduce the model complexity compared to the best model from **3.2**.

Depthwise separable convolution is a type of convolutional neural network (CNN) layer that consists of two distinct operations: a depthwise convolution and a pointwise convolution, which can reduce computational cost and memory usage.

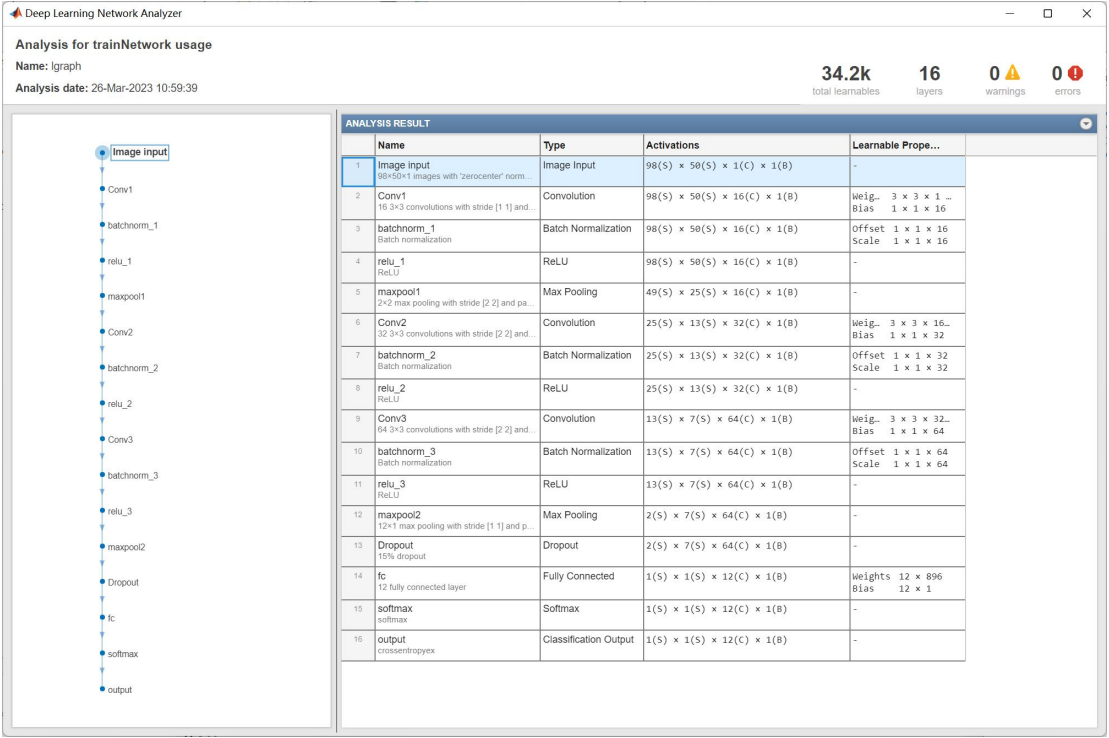
The first operation in a depthwise separable convolution layer is the depthwise convolution. This is a spatial convolution operation that applies a separate convolutional filter to each input channel, resulting in a set of output feature maps that are channel-wise independent. Essentially, it applies a filter to each channel of the input image independently, which allows the network to learn more specialized features for each channel. The depthwise convolution thus reduces the number of parameters in the network by a factor equal to the number of output channels.

The second operation in a depthwise separable convolution layer is the pointwise convolution. This is a 1x1 convolution operation that applies a linear transformation to the depthwise -convolved output, mapping it to a new feature space. This step can be thought of as a conventional layer, but applied to the output of the depthwise convolution instead of the original input image.

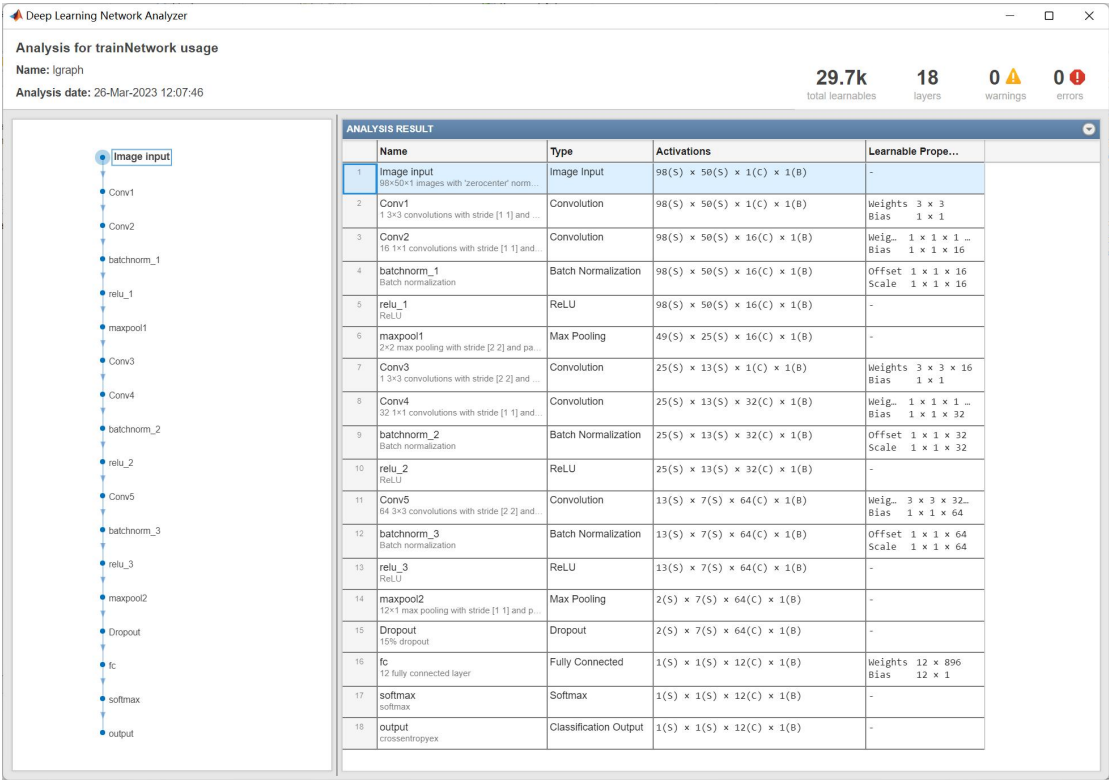
The combination of depthwise and pointwise convolutions enables a more efficient and effective learning process for CNNs. By separating the channel-wise and spatial information, depthwise

separable convolutions can learn more specialized features with fewer parameters than traditional convolutional layers, leading to faster training and better generalization performance.

Compare with CNNs with traditional convolutional layers:



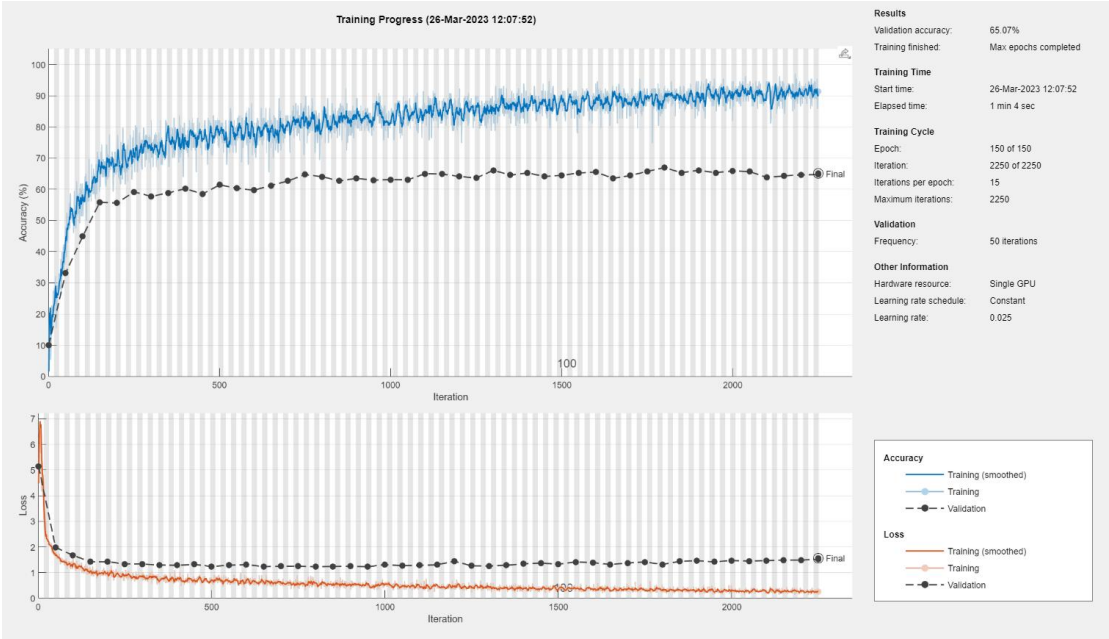
CNN with traditional convolutional layers



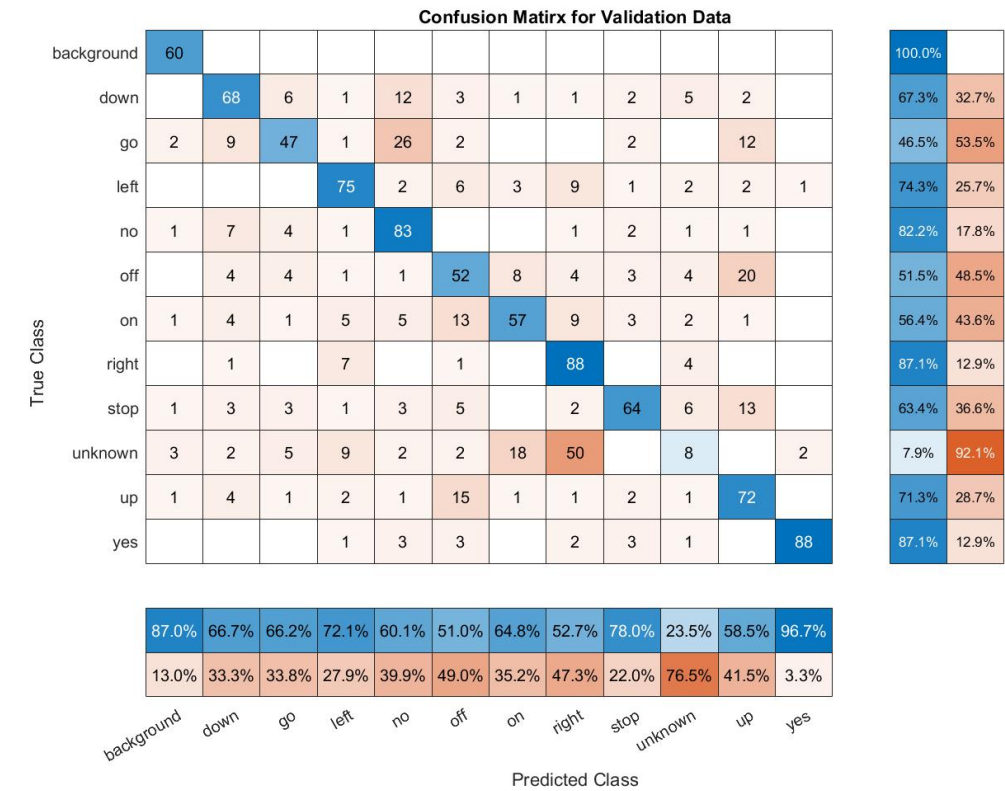
CNN with depthwise separable convolutional layers

As the experiment results shows, the CNN with traditional convolutional layers has 34.2k learnabels, whereas the CNN with depthwise separable convolutional layers has 29.7k learnables. Hence, the refined model has been reduced computational complexity comparing to the best model in 3.2.

Training progress:



Confusion plot:



4. Code:

Matlab:

```
1. clear all;
2. close all force;
3.
4. % define the random number seed for repeatable results
5. rng(1,'twister');
6. timePoolSize = 12;
7.
8. %% Load Speech Data
9.
10. % create an image data store from the raw images
11. imdsTrain = imageDatastore('C:\Users\Administrator\Desktop\DL_for_speech\speechImageData\speechImageData\TrainData',...
    'IncludeSubfolders',true,'LabelSource','foldernames')
12.
13.
14. % create an image validation data store from the validation images
15. imdsVal = imageDatastore('C:\Users\Administrator\Desktop\DL_for_speech\speechImageData\speechImageData\ValData',...
    'IncludeSubfolders',true,'LabelSource','foldernames')
16.
17.
18.
19.
20. % your code here...
21. % build network
22.
23. % 2x2 grid search:
24. % conv layers 2 + fliters 8 conv layers 3 + fliters 8
25. % conv layers 2 + fliters 16 conv layers 3 + fliters 16
26.
27. % The number of convolution layers: 2
28. % The number of filters of first convolution layer:8
29. % Net= [
30. %     imageInputLayer([98 50 1],'Name','Image input')
31. %     convolution2dLayer([3 3],8,'Padding','same','Name','Conv1')
32. %     batchNormalizationLayer('Name','batchnorm_1')
33. %     reluLayer('Name','relu_1')
34. %     maxPooling2dLayer(2,'Stride',2,'Padding','same','Name','maxpool1')
35. %
36. %     convolution2dLayer([3 3],16,'Padding','same','Stride', 2 , 'Name','Conv2')
37. %     batchNormalizationLayer('Name','batchnorm_2')
38. %     reluLayer('Name','relu_2')
39. %     maxPooling2dLayer([timePoolSize 1],'Name','maxpool2')
40. %
```

```

41. % dropoutLayer(0.15,'Name','Dropout')
42. % fullyConnectedLayer(12,'Name','fc')
43. % softmaxLayer('Name','softmax')
44. % classificationLayer('Name','output')
45. % ];
46.
47. % The number of convolution layers: 2
48. % The number of filters of first convolution layer:16
49. % Net= [
50. %     imageInputLayer([98 50 1],'Name','Image input')
51. %     convolution2dLayer([3 3],16,'Padding','same','Name','Conv1')
52. %     batchNormalizationLayer('Name','batchnorm_1')
53. %     reluLayer('Name','relu_1')
54. %     maxPooling2dLayer(2,'Stride',2,'Padding','same','Name','maxpool1')
55. %
56. %     convolution2dLayer([3 3],32,'Padding','same','Stride', 2 ,'Name','Conv2')
57. %     batchNormalizationLayer('Name','batchnorm_2')
58. %     reluLayer('Name','relu_2')
59. %     maxPooling2dLayer([timePoolSize 1],'Name','maxpool2')
60. %
61. %     dropoutLayer(0.15,'Name','Dropout')
62. %     fullyConnectedLayer(12,'Name','fc')
63. %     softmaxLayer('Name','softmax')
64. %     classificationLayer('Name','output')
65. % ];
66.
67. % The number of convolution layers: 3
68. % The number of filters of first convolution layer:8
69. % Net= [
70. %     imageInputLayer([98 50 1],'Name','Image input')
71. %     convolution2dLayer([3 3],8,'Padding','same','Name','Conv1')
72. %     batchNormalizationLayer('Name','batchnorm_1')
73. %     reluLayer('Name','relu_1')
74. %     maxPooling2dLayer(2,'Stride',2,'Padding','same','Name','maxpool1')
75. %
76. %     convolution2dLayer([3 3],16,'Padding','same','Stride', 2 ,'Name','Conv2')
77. %     batchNormalizationLayer('Name','batchnorm_2')
78. %     reluLayer('Name','relu_2')
79. %     convolution2dLayer([3 3],32,'Padding','same','Stride', 2 ,'Name','Conv3')
80. %     batchNormalizationLayer('Name','batchnorm_3')
81. %     reluLayer('Name','relu_3')
82. %     maxPooling2dLayer([timePoolSize 1],'Name','maxpool2')
83. %
84. %     dropoutLayer(0.15,'Name','Dropout')

```

```

85. % fullyConnectedLayer(12,'Name','fc')
86. % softmaxLayer('Name','softmax')
87. % classificationLayer('Name','output')
88. % ];
89.
90. % The number of convolution layers: 3
91. % The number of filters of first convolution layer:16
92. % Net= [
93. %     imageInputLayer([98 50 1],'Name','Image input')
94. %     convolution2dLayer([3 3],16,'Padding','same','Name','Conv1')
95. %     batchNormalizationLayer('Name','batchnorm_1')
96. %     reluLayer('Name','relu_1')
97. %     maxPooling2dLayer(2,'Stride',2,'Padding','same','Name','maxpool1')
98. %
99. %     convolution2dLayer([3 3],32,'Padding','same','Stride', 2 ,'Name','Conv2')
100. %     batchNormalizationLayer('Name','batchnorm_2')
101. %     reluLayer('Name','relu_2')
102. %     convolution2dLayer([3 3],64,'Padding','same','Stride', 2 ,'Name','Conv3')
103. %     batchNormalizationLayer('Name','batchnorm_3')
104. %     reluLayer('Name','relu_3')
105. %     maxPooling2dLayer([timePoolSize 1],'Name','maxpool2')
106. %
107. %     dropoutLayer(0.15,'Name','Dropout')
108. %     fullyConnectedLayer(12,'Name','fc')
109. %     softmaxLayer('Name','softmax')
110. %     classificationLayer('Name','output')
111. % ];
112.
113. % Reduce computational complexity by depthwise separable convolutions and
114. % 1x1 convolutions.
115. Net= [
116.     imageInputLayer([98 50 1],'Name','Image input')
117.     % depthwise separable convolutions
118.     convolution2dLayer([3 3],1,'Padding','same','Name','Conv1')
119.     convolution2dLayer([1 1],16,'Padding','same','Name','Conv2')
120.     batchNormalizationLayer('Name','batchnorm_1')
121.     reluLayer('Name','relu_1')
122.     maxPooling2dLayer(2,'Stride',2,'Padding','same','Name','maxpool1')
123.
124.     % depthwise separable convolutions
125.     convolution2dLayer([3 3],1,'Padding','same','Stride', 2 ,'Name','Conv3')
126.     convolution2dLayer([1 1],32,'Padding','same','Name','Conv4')
127.     batchNormalizationLayer('Name','batchnorm_2')
128.     reluLayer('Name','relu_2')

```

```

129.     convolution2dLayer([3 3],64,'Padding','same','Stride', 2 , 'Name','Conv5')
130.     batchNormalizationLayer('Name','batchnorm_3')
131.     reluLayer('Name','relu_3')
132.     maxPooling2dLayer([timePoolSize 1], 'Name','maxpool2')
133.
134.     dropoutLayer(0.15, 'Name','Dropout')
135.     fullyConnectedLayer(12, 'Name','fc')
136.     softmaxLayer('Name','softmax')
137.     classificationLayer('Name','output')
138. ];
139.
140. %Visualize the network
141. lgraph = layerGraph(Net);
142. analyzeNetwork(lgraph);
143.
144. % Configuration
145. options = trainingOptions('sgdm', ...
146.     'InitialLearnRate',0.025, ...
147.     'LearnRateDropFactor',0.1, ...
148.     'LearnRateDropPeriod',20, ...
149.     'MaxEpochs',150, ...
150.     'Shuffle','every-epoch', ...
151.     'ValidationData',imdsVal, ...
152.     'ValidationFrequency',50, ...
153.     'Verbose',true, ...
154.     'Plots','training-progress');
155.
156. net = trainNetwork(imdsTrain,Net,options);
157.
158. % evaluation
159. YPred = classify(net,imdsVal);
160. YValidation = imdsVal.Labels;
161. accuracy = sum(YPred == YValidation)/numel(YValidation);
162. fprintf('Accuracy: %.2f%%\n',accuracy*100);
163.
164. figure
165. cm = confusionchart(YValidation,YPred);
166. cm.ColumnSummary ="column-normalized";
167. cm.RowSummary = "row-normalized";
168. title('Confusion Matirx for Validation Data');
169. % cm.Normalization = 'absolute';

```