# Speech Commands Recognition

Xi Deng (xd487)*     Zulong Ye (zy2010)†

# 1  Introduction

## 1.1  Problem Description and Motivation

Nowadays, with the popularization of smart devices, speech recognition has become a promising way to reduce people's screen addiction. Speech recognition provides an effective approach for us to hands-free interact with all those screens, and the most important factor of voice recognition is its accuracy. It is never easy to transcript audio speech because human voices vary a lot from person to person in terms of accents, tones, paces and so on. In this project, we plan to develop an algorithm based on CNN or RNN with LSTM(Long Short-term Memory) that understands simple spoken commands. We will use the Speech Commands Dataset from Google TensorFlow, which includes 65,000 utterances by thousands of different people. We need to predict the exact command labels in test data "yes", "no", "up", "down", etc. after training.

## 1.2  Literature Survey

We have reviewed three paper performing experiments on the TIMIT phoneme recognition dataset. In *Convolution Neural Networks for Speech Recognition*, the authors show us how to implement CNN on speech input and they manage to reach the average error rate of 20%. While the authors achieve a lower error rate of 17.7% by applying deep bidirectional RNN with LSTM in *Speech Recognition with Deep Recurrent Neural Networks*. Both have noticeable performance improvements over traditional DNNs on the same dataset. In *Research and System Design of Speech Recognition Based on Improved CNN*, another author has done the comparison between CNN with GPU acceleration and RNN with LSTM. As a result, the improved CNN can not only save the running time by 20%, but also reduce the recognition error rate by 15%.

*Tandon School of Engineering, New York University. Email: xd487@nyu.edu.

†Tandon School of Engineering, New York University. Email: zy2010@nyu.edu.

# 2   Technical Details

## 2.1   Datasets

The Speech Commands Dataset, released by Google TensorFlow, is a set of one-second wav audio files, which contains 65,000 samples of human voices. In this dataset, files are organized into folders based on the words spoken (labels). There are 30 kinds of short words including core command words such as "yes", "no", "up", "down", "left", "right", "on", "off", "stop", "go", from "zero" to"nine" and auxiliary words such as "bed", "bird", "cat", "dog", "happy", etc. Since the size of the dataset is very large, we plan to extract a part of the data as an experiment to test whether our model works. To validate our model's performance, we will use 80% of the whole dataset for training and 20% of it for testing.

In terms of these .wav file speech information, we use 2 methods to read them. Firstly, we use wavfile.read() function to get the spectrogram with 16000Hz and normalize it. As it takes too long a time to get all data, we store the data information to a .txt file which includes all .wav files information. Then, we can use numpy arrays to read the .txt file to get all the data with this second method, which improves the reading speed greatly.

## 2.2   Models

**Input**: For wav audio files, we can use the spectrogram to analyze them. All files can be extracted to a numpy array with a sample rate of 16000Hz. Then, it can be regarded as 16000 * 1 input features for each sample.

**Model implementation**: In the convolution neural network, we can set the input dimension to be 16000 * 1. Next, we can construct a few convolution 1D layers and some fully connected layers. In the convolution layer, we can take "same" padding mode and "relu" as the activation function. In the fully connected layer, we will also take "relu" as the activation function. In the output layer, "softmax" is the activation function.

In the LSTM model, we can also set the input dimension of 16000 * 1. We implement 2 layers of LSTM. All parameters need to be adjusted for better accuracy.

**Output**: The output result is the probability of each class representative, and we will select the largest probabilities as the predicted label. For actual values computed by the label of each file folder, we will match each text label with a number. Then, each label can be matched with a class representative.

## 2.3   Algorithms

Firstly, we denoise the speech data using a high pass filter and resample the input so that it can be taken by CNN and RNN. Then we train and modify the model so that it can make the right prediction of ten commands "yes", "no", etc., with a high accuracy. Finally, we apply both models to test data and compare the results of both networks (pure CNN and CNN with LSTM). The main flow chart diagram of our project is shown in Figure 1.
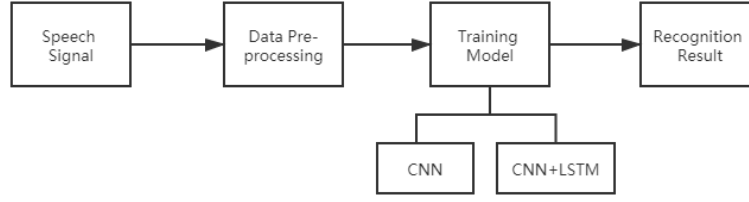
FIGURE 1: Flow chart diagram

In terms of data pre-processing, we choose "yes", "no", "up", "down", "left", "right", "on", "off", "stop", "go" commands as the target for recognition. By using the "LabelEncoder" to change each text label into specific numbers, we regard the value of 1-10 numbers as ten classes representative.

Since we aim to predict the probability of the label, we would choose "cross entropy" as the loss function, and "stochastic gradient descent" as an optimizer for CNN; "cross entropy" as the loss function and "adam" as an optimizer for LSTM. Through the loss function and "softmax" algorithm, we can calculate the probability for each speech command, and choose the largest probability as the predictive value.

For the convolutional 1D layer, we set the 'Relu' as the activation function. Therefore, we have all neuron activations in each layer can be represented in the following matrix form:

$$o^{(l)} = \phi(o^{(l-1)}W^{(l)})(l = 1, 2, ..., L-1),$$

where $\phi(x) = max(0, x)$. The loss function shows as follows.

$$L_{CE} = -\sum_{i=1}^{n} t_i log(p_i),$$

for n classes, where $t_i$ is the actual label and $p_i$ is the Softmax probability for the $i^{th}$ class. And the softmax probability is given by:

$$p_i = \frac{exp(y_i)}{\sum_{j=1}^{n} exp(y_j)}$$

## 2.4   Training Methods

**Training Models (CNN):**
We finalize our CNN model including 10 convolution 1D layers and 1 fully connected layers to achieve the best training result. We set our parameters as follows. From the first to the ninth layer, there are 8, 16, 32. . . 2048 hidden units and a 9*9-size filter for each layer. Besides, we apply Batch Normalization to these layers, and we set the max pooling size to 2*2 for each layer. When it comes to the tenth convolution layer, we use 1024 hidden layers and there is a fully connected dense layer with 1024 units. We choose 'same' padding mode

FIGURE 2: CNN model and training details



FIGURE 3: CNN+LSTM model and training details

for every conv1D layer, and we use dropout with 0.2 to reduce overfit. In addition, a dense layer with 10 units is added for label classifications in output. Last but not least, we plan to train the model for 300 epochs with "early stopping" callback function. Model and training details are shown in Figure 2.

**Training models (CNN + LSTM):**
We finalize our CNN and LSTM model with 1 convolution 1D layer with 256 hidden units and a 32*32-size filter using the "same" padding mode, as well as 2 LSTM layers with 256 units connected to the first convolution layer. Next, there is 1 fully connected layer with 64 units. We use dropout with 0.2 to reduce overfit. Lastly, a dense layer with 10 units is added for label classifications. Model and training details are shown in Figure 3.

## 2.5   Evaluation Metrics

We are using training accuracy, validation accuracy, testing accuracy and confusion matrix to evaluate the performance of our CNN and LSTM models. The index of confusion matrix is given by the number of correct predicted labels divided by the number of actual labels. The vertical line is the predicted labels and horizontal line is the actual label.

# 3   Results

## 3.1   Data Visualization

We can open a speech command file through the "wavfile.read()" function and convert it into a numpy array. To extract its features, we plot its wave spectrum and spectrogram using the "log_spectgram()" function. For example, we take a .wav file from "left" commands, then plot its wave spectrum and spectrogram in Figure 4.



FIGURE 4: Wave spectrum and spectrogram of a speech command

Next, we will count the number of audio files in each label, as shown in the left part of Figure 5. Then, we split the whole dataset as training data, validation data, and test data with a ratio of 0.2 test size. The right part of Figure 5 represents the result of data split.



```
training data shape: (12786, 8000, 1)
validation data shape: (4263, 8000, 1)
test data shape: (4263, 8000, 1)
```

(a) Each label                              (b) Data size

FIGURE 5: Number of samples in each label and training, validation, testing data size

## 3.2    Train Results

**Training Models (CNN):**

After a number of attempts on adjusting the CNN model structure and coefficients, we manage to reach a training accuracy of 99%, a validation accuracy of 94% and a test accuracy of 93%. The curve of training accuracy is smoother and converges at about 15 epochs, while the curve of validation accuracy bounces a lot. From the confusion matrix on the right, we can see that the accuracy of predicting each label is higher than 90%, while the highest one goes above 96%.



```
[ ]  model_cnn.evaluate(x=x_test, y=y_test)

     134/134 [==============================] - 2s 16ms/step - loss: 0.2937 - accuracy: 0.9353
     [0.2936517894268036, 0.935256838798523]
```

(a)  Training result



(b)  Confusion matrix

FIGURE 6:  CNN training results

**Training Models (CNN+LSTM):**

We obtain a training accuracy of 96%, a validation accuracy of 90%, and a test accuracy

of 91%. The curve of both training and validation accuracy are smooth, but validation accuracy seems to converge earlier than training accuracy. From the confusion matrix on the right, we can tell that the prediction accuracy of each label is around 90%, and the highest one is 96% while the lowest one is 88%.



```
[ ]  model_lstm.evaluate(x_test, y_test)

     134/134 [==============================] - 2s 12ms/step - loss: 0.4494 - accuracy: 0.9113
     [0.4494383633136749, 0.9113300442695618]
```

(a) Training result



(b) Confusion matrix

FIGURE 7: CNN+LSTM training results

The following table demonstrates the comparison between the performance of some models that we have tested and the performance of our finalized model in CNN and CNN+LSTM.

TABLE 1: Training Results of various network structures

| | Network Structure | Average Test Accuracy | Number of Training Parameters | Average Running Time |
|---|---|---|---|---|
| 1 | CNN (2 conv1D, 1 dense, input size = 8000) | 59.39% | 7.351M | 8m 43s |
| 2 | CNN (4 conv1D, 2 dense, input size = 16000) | 78.14% | 1.0996M | 22m 59s |
| 3 | CNN (10 conv1D, 2 dense, input size = 16000) | 93.5% | 28.336M | 45m 21s |
| 4 | CNN+LSTM (1 conv1D, 1 LSTM, 1 dense, input size = 8000) | 86.61% | 14.676M | 3m 54s |
| 5 | CNN+LSTM (1 conv1D, 2 LSTM, 1 dense, input size = 16000) | 90.17% | 3.108M | 3m 37s |
| 6 | CNN+LSTM (1 conv1D, 4 LSTM, 1 dense, input size = 16000) | 89.96% | 4.158M | 5m 27s |

## 3.3 Predict

Here is a list of predictions of a couple of speech command samples using both our models. We can see that both our models make predictions very precisely, producing almost no mistake.

```
Audio: go            Audio: on            Audio: off           Audio: stop
Text (CNN): go       Text (CNN): on       Text (CNN): off      Text (CNN): stop
Text (LSTM): go      Text (LSTM): on      Text (LSTM): off     Text (LSTM): stop
Audio: on            Audio: no            Audio: go            Audio: up
Text (CNN): on       Text (CNN): no       Text (CNN): go       Text (CNN): up
Text (LSTM): on      Text (LSTM): no      Text (LSTM): go      Text (LSTM): up
Audio: on            Audio: left          Audio: no            Audio: stop
Text (CNN): on       Text (CNN): left     Text (CNN): no       Text (CNN): up
Text (LSTM): on      Text (LSTM): left    Text (LSTM): no      Text (LSTM): right
Audio: up            Audio: up            Audio: left          Audio: up
Text (CNN): up       Text (CNN): up       Text (CNN): left     Text (CNN): up
Text (LSTM): up      Text (LSTM): up      Text (LSTM): stop    Text (LSTM): up
Audio: down          Audio: down          Audio: stop          Audio: go
Text (CNN): down     Text (CNN): down     Text (CNN): stop     Text (CNN): go
Text (LSTM): down    Text (LSTM): down    Text (LSTM): stop    Text (LSTM): go
Audio: left          Audio: yes           Audio: no            Audio: up
Text (CNN): left     Text (CNN): yes      Text (CNN): no       Text (CNN): up
Text (LSTM): left    Text (LSTM): yes     Text (LSTM): no      Text (LSTM): off
Audio: yes           Audio: go            Audio: go            Audio: off
Text (CNN): yes      Text (CNN): go       Text (CNN): go       Text (CNN): off
Text (LSTM): yes     Text (LSTM): go      Text (LSTM): go      Text (LSTM): off
Audio: go            Audio: on            Audio: yes           Audio: no
Text (CNN): go       Text (CNN): on       Text (CNN): yes      Text (CNN): no
Text (LSTM): go      Text (LSTM): on      Text (LSTM): yes     Text (LSTM): no
Audio: go            Audio: up            Audio: yes           Audio: up
Text (CNN): go       Text (CNN): up       Text (CNN): yes      Text (CNN): up
Text (LSTM): go      Text (LSTM): up      Text (LSTM): yes     Text (LSTM): up
Audio: no            Audio: stop          Audio: no            Audio: go
Text (CNN): no       Text (CNN): stop     Text (CNN): no       Text (CNN): go
Text (LSTM): no      Text (LSTM): stop    Text (LSTM): no      Text (LSTM): go
Audio: go            Audio: no            Audio: down          Audio: off
Text (CNN): stop     Text (CNN): no       Text (CNN): down     Text (CNN): off
Text (LSTM): off     Text (LSTM): no      Text (LSTM): down    Text (LSTM): off
Audio: yes           Audio: no            Audio: yes           Audio: right
Text (CNN): yes      Text (CNN): no       Text (CNN): yes      Text (CNN): right
Text (LSTM): yes     Text (LSTM): no      Text (LSTM): yes     Text (LSTM): right
Audio: right         Audio: stop          Audio: no            Audio: off
Text (CNN): right    Text (CNN): stop     Text (CNN): no       Text (CNN): up
Text (LSTM): right   Text (LSTM): stop    Text (LSTM): no      Text (LSTM): up
Audio: down          Audio: right         Audio: stop          Audio: stop
Text (CNN): down     Text (CNN): right    Text (CNN): stop     Text (CNN): stop
Text (LSTM): down    Text (LSTM): right   Text (LSTM): off     Text (LSTM): stop
Audio: right         Audio: on            Audio: left          Audio: on
Text (CNN): right    Text (CNN): on       Text (CNN): left     Text (CNN): on
Text (LSTM): down    Text (LSTM): on      Text (LSTM): left    Text (LSTM): on
```

FIGURE 8: Predict results

# 4 Conclusions and Source Code

## 4.1 Conclusions

In this project, we manage to implement a deep convolutional neural network and a combination of CNN with LSTM network, as well as apply both networks on the Speech Commands Dataset from Google TensorFlow. Now we have accomplished the feature extractions of audio files, and increased the speed of audio data processing greatly. Most importantly, we have acquired very high recognition precisions on both finalized networks, which shows a significant improvement on what we did in the project update. As a result, the complex CNN model has slightly better performance compared to the CNN with LSTM model, while the CNN with LSTM model runs much faster than the complex CNN model. Furthermore, we found that CNN performance is sensitive to padding mode and pooling size, rather than filter size. In addition, CNN performance can be enhanced by adding more and more layers to make the network very deep but LSTM cannot.

## 4.2 Future Work

We plan to apply our system on some large-vocabulary speech dataset, or try some other feature extraction approaches like Mel Frequency Cepstrum Coefficient(MFCC) in the future. We may also try some new models such as Transformers/BERT, and we can compare the performance of these different models.

## 4.3 Source Code

Github: https://github.com/yzleaf/speech_commands_recognition
Colab (Access by NYU account): Click here

# References

1. O. Abdel-Hamid, A. Mohamed, H. Jiang, L. Deng, G. Penn and D. Yu, "Convolutional Neural Networks for Speech Recognition," in IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 22, no. 10, pp. 1533-1545, Oct. 2014, doi: 10.1109/TASLP.2014.2339736.

2. A. Graves, A. Mohamed and G. Hinton, "Speech recognition with deep recurrent neural networks," 2013 IEEE International Conference on Acoustics, Speech and Signal Processing, Vancouver, BC, Canada, 2013, pp. 6645-6649, doi: 10.1109 /ICASSP.2013.6638947.

3. L. Xu, "Research and System Design of Speech Recognition Based on Improved CNN," Master, Donghua University, Shanghai, China, 2019.

4. https://www.tensorflow.org/datasets/catalog/speech_commands

5. https://keras.io/api/layers/recurrent_layers/lstm/

6. Davids, "Speech representation and data exploration", https://www.kaggle.com/davids1992/speech-representation-and-data-exploration (accessed May 2021)