

## Abstract

In the past a few months, I and David were together working on a series of tasks, including simple Audio to Text conversion, software tests on android device, Cloud Infrastructure, Kaldi platform setup, increasing training data, neural network pruning. We also spent time on Iflytek SDK development, but it will not to be discussed in this report. At current stage, we are facing challenge in increasing data and dropping connection.

### **Task1: Simple Audio to Text conversion**

Our first task is to train a model to recognize 10 Chinese words from the Chinese class vocabulary using TensorFlow. The idea of training audio processing model with CNN architecture comes from the revolutionary keyword spotting method of Google Inc. New York.<sup>[1]</sup> With this method, we trained a CNN model using spectrogram<sup>[2]</sup>. We trained 105,000 WAVE audio datasets through 18,000 epochs. The final accuracy is 88.2%(N=4890). The outcome confusion matrix is given below (Fig1.1). We test this model on android device. It is proved to be able to recognize 10 English words with a relatively high accuracy.

```

INFO:tensorflow:Step #18000: rate 0.000100, accuracy 80.0%, cross entropy 0.494432
INFO:tensorflow:Confusion Matrix:
[[366  1  1  0  1  0  0  0  0  0  1  1  0]
 [ 2 275  6  5 12 12  8 15 10  3 12 11]
 [ 1  3 370  9  1  3  8  1  0  0  0  1]
 [ 4  9  4 347  4  7  6  0  0  1  5 19]
 [ 3  3  0  0 320  1  3  1  3  6  8  2]
 [ 2  3  1 19  1 335  1  0  0  1  3 11]
 [ 2  2 13  3  2  0 320  6  0  2  1 1]
 [ 1  8  2  0  2  0  8 339  1  1  0 1]
 [ 5 10  0  0  8  1  0  0 335  4  0 0]
 [ 4  3  0  0 23  0  3  2  9 323  4 2]
 [ 5  3  0  1 13  0  1  1  0  2 320 4]
 [ 2 16  0 29  4 12  0  1  2  1  3 302]]
INFO:tensorflow:Step 18000: Validation accuracy = 88.9% (N=4445)
INFO:tensorflow:Saving to "/tmp/speech_commands_train/conv.ckpt-18000"
INFO:tensorflow:set_size=4890
INFO:tensorflow:Confusion Matrix:
[[405  2  0  0  0  0  0  0  0  0  0  1  0]
 [ 3 303  6  8  7 11 17 13  8  1 11 20]
 [ 0  9 394  2  0  1  9  0  0  0  1  3]
 [ 0  6  7 344  2 15  7  2  0  0  1 21]
 [ 2  5  0  1 390  4  2  0  4  7  9 1]
 [ 1 11  4 30  1 346  3  0  1  0  1 8]
 [ 1  3 17  0  3  0 381  5  0  0  2 0]
 [ 1 11  1  1  0  0  6 370  4  2  0 0]
 [ 0  7  0  0  4  9  0  0 353 16  2 5]
 [ 1  2  0  0 27  0  4  1 12 346  4 5]
 [ 0  3  0  1 12  5  0  0  0  3 382 5]
 [ 3 15  0 56  4 12  3  4  0  1  3 301]]
INFO:tensorflow:Final test accuracy = 88.2% (N=4890)
yichuanj95@instance-tensorflow2:~/tensorflow$ 
[0] 0:bach* "instance-tensorflow2" 05:3

```

Fig1.1 outcome confusion matrix

## Task2: Cloud Infrastructure

Since the training process takes days to complete, we use Google Cloud's GPU to finish the long-time training process. Also, it's more secure to store our model on the cloud. We store our model in different instances across multiple private zones. Tmux is used to master the session in the training process. We keep snapshots for models.

## Task3: Network Pruning and connection dropping

We tried to do weight pruning with Keras API. But very soon we realized that Keras weight pruning API only targets MNIST model<sup>[3]</sup>. However, our trained model is MobileNetV1. It means we can not use it to do compression. Also, we are now struggling how to set

the `initial_sparsity`, `begin_step`, `end_step` and frequency in the `pruning_schedule`.<sup>[4]</sup>

<https://stackoverflow.com/questions/56470335/is-there-any-keras-code-to-reproduce-the-weight-pruning-of-mobilenet>

#### **Task4: Increasing Data**

This section focused on making artificial voice data. Since voice data is deficient, generating artificial data becomes a necessary prerequisite for further model training. Using the code from [here](#) We are still struggling with debugging the code. The code is based on python 2.7 which will be expired by 2020. Some of the libraries are not compatible in python3.6 and windows environment, so we tried two solutions. The first one is to use another code which is for python3, the second one is to use the linux environment and python2.7 on the virtual machine.

## Notations:

[1]: [https://www.tensorflow.org/tutorials/sequences/audio\\_recognition](https://www.tensorflow.org/tutorials/sequences/audio_recognition)

[2]: Convolutional Neural Networks for Small-footprint Keyword Spotting, ara N. Sainath, Carolina Parada Google, Inc. New York, NY, U.S.A

[3]: <https://github.com/tensorflow/tensorflow/issues/29163>

[4]: <https://stackoverflow.com/questions/56470335/is-there-any-keras-code-to-reproduce-the-weight-pruning-of-mobilenet>

## References:

1. Convolutional Neural Networks for Small-footprint Keyword Spotting, ara N. Sainath, Carolina Parada Google, Inc. New York, NY, U.S.A
2. J. Schalkwyk, D. Beeferman, F. Beaufays, B. Byrne, C. Chelba, M. Cohen, M. Kamvar, and B. Strobe, ““Your word is my command”: Google search by voice: A case study,” in Advances in Speech Recognition, A. Neustein, Ed. Springer US, 2010, pp. 61–90.
3. G. Chen, C. Parada, and G. Heigold, “Small-footprint Keyword Spotting using Deep Neural Networks,” in Proc. ICASSP, 2014