# Speech Diarization using Deep Neural Networks and Machine Learning

Clint Ferrin, Madi Mickelsen, Chui Vanfleet, Daniel Mortensen

December 13, 2017

# Contents

# List of Figures

# 1  Introduction

"Audio diarization is the process of annotating an input audio channel with information that attributes (possibly overlapping) temporal regions of signal energy to their specific sources" [1]. This report represents an effort to solve speaker diarization problem using deep learning. Many other methods of speech diarization rely heavily one signal processing techniques, but fewer attempt to use neural networks [2, 3, 4].

Dr. Stephanie A. Borrie from the Human Interaction Lab at Utah State University has provided a number of audio files from her research. The audio files contain a recording of two individuals conversing. Each individual speaks into their own microphone. But, crosstalk from each speaker is picked up on the other microphone.

Along with the audio files, Dr. Borrie has provided TextGrid files which contain the diarization annotations. These TextGrid files were created by Dr. Borrie's students and are to be used as training and testing data for a deep learning algorithm.

## 1.1  Problem

Using the audio and TextGrid files provided by Dr. Borrie, a deep learning algorithm must be developed to automatically perform speaker diarization

# 2  Requirements

## 2.1  Algorithm Requirements

1. The project must use a neural network at some point. If additional processing steps are used, such as in preparation of the data, or organizing the output response, that is also fine.

2. The training code should provide some indication of how the training is progressing.

3. The code should write a file indicating the diarization.

## 2.2  Report Requirements

1. Clearly describe the architecture used in processing.

2. Clearly describe the neural architecture, such as number of layers, neurons in each layer, type of neuron, number of weights, stride, etc.

3. Clearly describe what the inputs are.

4. Clearly describe how many outputs there are, and what the output layer function is (sigmoidal, softmax, linear, etc.).

5. Clearly describe the cost function the neuron is training against, and the kind of optimization used.

6. Clearly describe any other processing steps used in conjunction with the neural net.

7. Describe how to train the neural net — how it is presented input data, how the output data is presented.

8. Present data, preferably learning curves, demonstrating the neural network is learning. Indicate how many training iterations used, and why.

9. Present results quantifying how well the diarizer works. Ideally, this would be compared to the original training data. Provide a description of how the score is derived.

10. Provide a description of what was learned.

11. Provide a printout of the code.

# 3 Methods

Several different diarization methods were explored to meet the requirements discussed in subsection 2.1. This section will discuss, in detail, the several diarization methods created, presented in no particular order.

## 3.1 Fully Connected Neural Network

The fully connected version of the Neural Network can be seen in Figure 1. The Network has 5 layers, it receives a total of 3400 inputs (1600 inputs from channel 1, 1600 inputs from channel 2, 100 inputs from the `fft` of channel 1, and 100 inputs from the `fft` of channel 2), and has four possible outputs (00-no speech, 01-speaker one, 10-speaker two, and 11-both speakers).

The input layer is stacked horizontally with channel 1 sampled at 8000 Hertz for a total length of 0.2 seconds, channel 2 sampled at 8000 Hertz, and finally both channels are followed by their corresponding `fft` results with a size of 100 each. This yields a total input size of 3600; each new sample of training or testing data is stacked horizontally on top off each input sample to be fed into the network—the stacking of input data in matrix form allows TensorFlow to speed up the training process.

In order from left to right, each layer has 300, 200, 200, and 500 neurons. All of the hidden layers utilize ReLU activation functions, and the densely connected layers 2 through 4 have a dropout rate of 0.8. The output layer uses the softmax activation function because there are 4 different classes that the loss function needs to categorize.

To train the network, a file path is given to the main function. The number of `.wav` files in the audio files directory must match the number of `.TextGrid` files in the labels
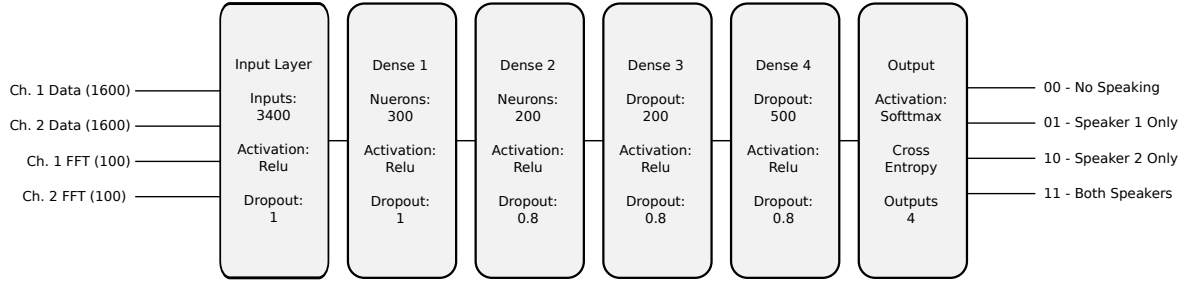
Figure 1: Block Diagram of Neural Network

directory. The label `numpy` array is created by sampling the `TextGrid` file to create an output list that corresponds to every $Fs$ (samples per second) in the audio file.

The cost function utilized in the network design is cross-entropy, and the output of the network is compared to a one-hot-encoded `Y` vector that contains the actual class that was marked by humans. The optimization function that performed best with this setup was the `AdamOptimizer` that is included in the TensorFlow parameters. For more details reference the TensorBoard printout in Figure 2
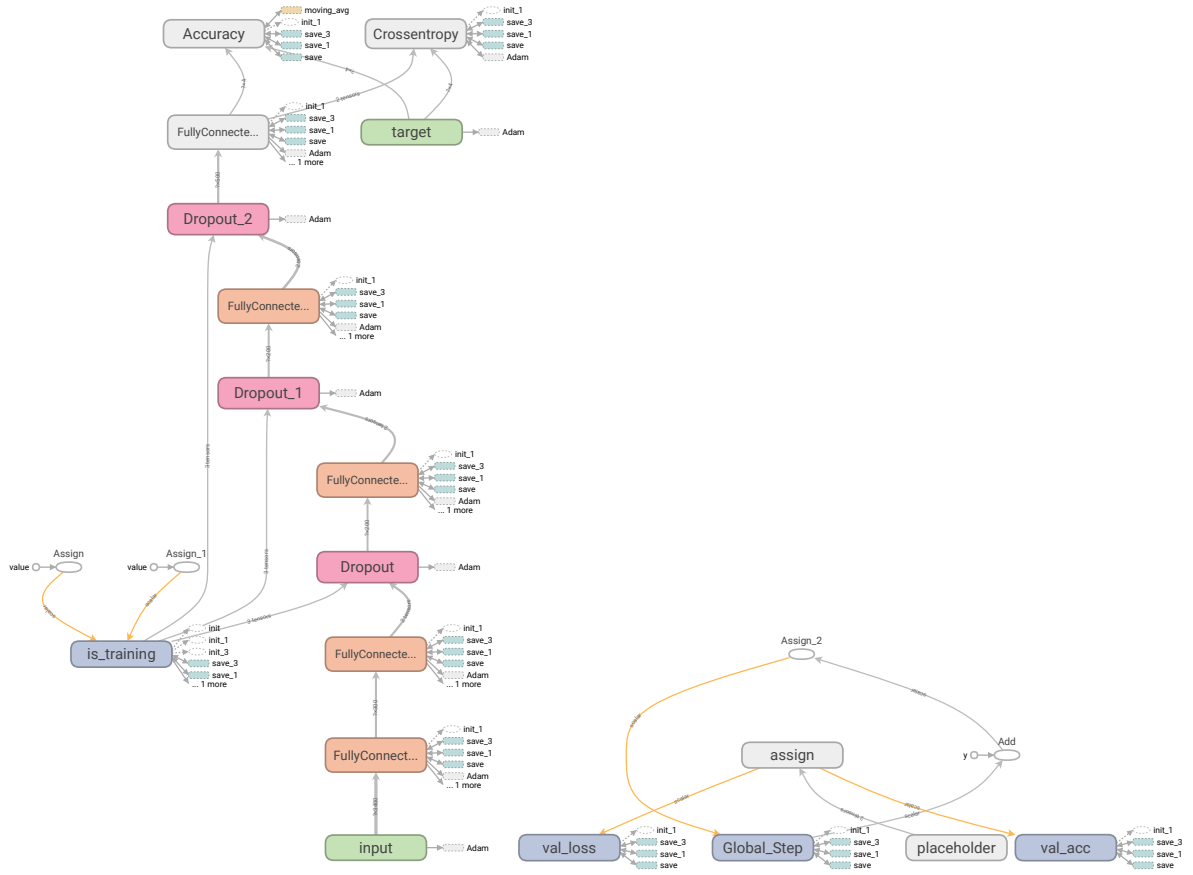


Figure 2: Fully Connected Neural Network

While implementing the network, different numbers of neurons, dropout rates, and

network types were experimented with. The 1D convolutional network did not produce better results, and if the dropout rate was too high, the loss function would not converge. The parameters were changed and re-run to identify the best network system.

Introducing the FFT significantly improved the accuracy of the network. Before the FFT, the network had a very difficult time identifying silence, but the added information from the FFT improved the overall performance. For the results of this method, see Section 4

## 3.2  Blind Source Separation

The original design included a preprocessing portion using blind source separation (BSS). The goal of BSS, also known as independent component analysis, is to determine a matrix which forces the elements of the output to be statistically independent. The matrix is determined using an iterative process of maximum likelihood independent component analysis. When multiplied by the output, each element should be separated from the others.

For the purpose of this project, BSS would be used to separate the two speakers when speech overlap occurs, removing crosstalk. The BSS program was developed for another class using MATLAB. Before implementing BSS into the design fully (converting the MATLAB program into python), segments of an audio file with overlapped speech were run through the BSS program developed in MATLAB. This would test the BSS method to determine if it would add, or detract from the system.

Initially, the program would not separate the signals whatsoever. The audio appeared to be of worse quality, with more crosstalk rather than less. More iterations were run, the step size was changed, and more segments were concatenated to see if more data worked, but none of these possible *improvements* did not appear to change anything.

The reason BSS did not work for this instance could be due to the delay between microphones. For example, speaker one speaks directly into microphone one, whereas speaker two speaks directly into microphone two. The distance between the speakers would cause a delay to occur between speaking into one microphone and the crosstalk on the second. The BSS algorithm used in the program does not account for the delay. As a result this method was not incorporated into the final design of the system.

## 3.3  Transfer Function Estimation

Based on previous results, the following is a writeup of the next step that is predicted to create better results then the previous. In past attempts, windows of data were inserted into a neural network, along with their corresponding FFTs, to determine where changes in speaking would occur. This was complicated by microphone recording speech signals from several participants at once. This section goes over a method designed to show which microphone is being used during a given time interval.

### 3.3.1 Methodology

This method operates under the assumption that the input to a single microphone consists of the input from the nearest speaker and the output from the second speaker after it has been run through a transfer function $H(Z)$ as seen in Figure 3. This portion
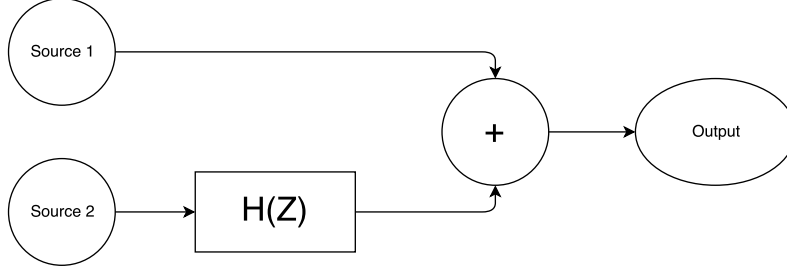


Figure 3: Overall System Diagram

of the speech diarization would be purposed towards detecting who is speaking: person one, person two, or both. For intents of the following setup, it is assumed that one person is speaking at a time. This algorithm requires two channels, one for microphone one, and the other for microphone two. In this setup, if only one person is speaking, then one signal will represent the most accurate data from the source, and the other will represent the same data after it has passed through a system and been modified by the system's transfer function.

The proposed solution is to develop a neural network that would simulate the transfer function that operates on the non-dominant dataset. Once this is accomplished, channel A is run through the simulated transfer function and subtracted from the other channel. If the output is minimal, then the channel that was operated on matches its user. In other words, when person one's data is run through the transfer function, and subtracted from person two's data, if the output is close to zero, then person two was not talking, and person one provided input to the microphone and vice versa. Figure 4 demonstrates a system view if the input signal was pre-transformed. Figure 5 shows the
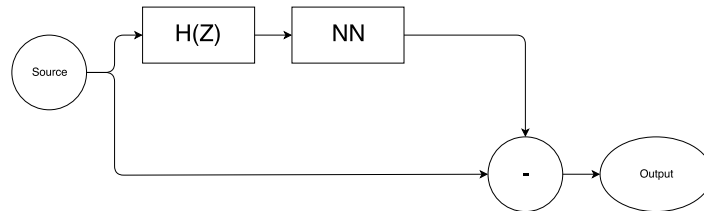


Figure 4: Non-Zero Error System

system view if the input signal is not passed through the transfer function, for example, it is person one's signal when person one was speaking. Channel two should then be run on input data. The results from both systems should be compared. If non-zero values

7

Figure 5: Zero Error System

are detected after running both channels through the transfer function, then both were talking. An example of the proposed system is shown in Figure 6. Both the original sig-



Figure 6: Neural Network System Design

.

nal is provided as input ,and the time delayed versions of the same signal. It is assumed when the signal undergoes a transformation, the transformed signal is consists of linear combinations of time delayed source samples.

### 3.3.2 System Parameters

The neural network will be comprised of a semi-connected layer where each neuron has access to a time delayed sample of the input data. The data is arranged as a Toeplitz matrix, $A$, with each row representing the input data of each neuron. The output of the network has equivalent dimensions as the output of the transfer function. There should be between 0 - 500 sample delays based on cross correlation, as shown in Figure 7. Stronger correlations between signals are seen from as far as $\tau = 500$ samples out.

8

Figure 7: Cross Correlation Data

# 4 Results

Using the network described in Section 3.1, the network correctly labeled a file it had not seen with 75% accuracy. The training produced the following loss function graph after 25 epochs. 25 epochs was chosen, with batch sizes of 100 because larger batches or epochs cause the network start to decline in accuracy due to over-fitting. The loss plot in Figure 8 converges to a value; with more epochs, little progress is made in reducing the loss.



Figure 8: TensorBoard Loss Training Plot

The output of the accuracy plot can be seen in Figure 9. The accuracy was obtained by splitting the data into random training and testing sets of length 0.2 seconds. To create the training and testing data sets, all of the data was placed into a long array of data and shuffled in the same way as the labels. The data was down-sampled to a new rate of 8000 Hz, and the labels were sampled at 5 Hz. After the shuffling, the data was split into the proportion indicated by the variable `percent_train`.

As a result, the accuracy represents categorizing data that has not been trained by

the Neural Network.



Figure 9: TensorBoard Accuracy Training Plot

It should be noted that the human markings were incorrect in many of the text grids. For example, the file `HS_D09.wav` seen in Figure 10 shows the first 10 seconds of a file that has been poorly marked. When listening to the file with both channels, it is obvious that no one is speaking for the majority of the 10 seconds. In that audio portions that people are speaking, the markings are incorrect. If the TextGrid files were more accurately labeled, the output of the programs would be more accurate as well.



Figure 10: Example of Poor Diarization

To test the effectiveness of the training, the network was given a complete file of

4:00 minutes, and told to predict the labels. Despite the poorly marked training data, the network achieves 75% accuracy on an entire data file that the network had not seen previously. To test the accuracy, the diarization state was predicted every 0.2 seconds and compared to the markings on each tier of the TextGrid files. As a visual representation, the output after the network predicted the labels of the file was graphed. Figure 11 shows all of the necessary information to view how the network performed.



Figure 11: Results of Network on Untrained Data

The solid fill colors represent where the human labeled the speaking of both people. The areas with no color represents when neither person one or person two are speaking, blue fill represents when person one is speaking, red fill represents when person two is speaking, and green fill represents when both people are speaking at the same time.

The solid lines at the bottom of the graph are the predictions that the network produces. The predictions of the network line up closely to the human markings. In many cases, the neural network identifies small sounds that people make that are not recorded as speech in the TextGrid file (such as brief laughs). This can be filtered out at a future date to improve accuracy.

# 5   Conclusion

The fully connected neural network with the fft inputs was used to label audio and performed with 75% accuracy. The fully connected neural network was chosen as our

final design because it yielded the best results. TensorFlow was a challenge to use by itself and it was difficult to quickly change network design, so using APIs like TFLearn and Keras made the process much easier.

Alternate methods to the speech diarization problem were discussed once the initial design was implemented and tested. These would be implemented in future iterations of this project to determine if they improved the accuracy.

The first alternate approach would use the Neural Network for the BSS. The NN should be able to incorporate the delay, making it more robust then the BSS algorithm used initially. This would allow easy separation of the two channels. After filtering out the noise, the speech on each channel would be easily detectable using the power of the signal.

Section 3.3 discussed an another method of approaching the problem, estimating the transfer function. The transfer function estimation method would help to identify when each microphone is being used but may prove difficult to implement depending on the different recording environments.

# References

[1] S. E. Tranter and D. A. Reynolds. "An overview of automatic speaker diarization systems". In: *IEEE Transactions on Audio, Speech, and Language Processing* 14.5 (Sept. 2006), pp. 1557–1565. ISSN: 1558-7916. DOI: 10.1109/TASL.2006.878256.

[2] K. Boakye et al. "Overlapped speech detection for improved speaker diarization in multiparty meetings". In: *2008 IEEE International Conference on Acoustics, Speech and Signal Processing*. Mar. 2008, pp. 4353–4356. DOI: 10.1109/ICASSP.2008.4518619.

[3] X. Anguera et al. "Speaker Diarization: A Review of Recent Research". In: *IEEE Transactions on Audio, Speech, and Language Processing* 20.2 (Feb. 2012), pp. 356–370. ISSN: 1558-7916. DOI: 10.1109/TASL.2011.2125954.

[4] C. Barras et al. "Multistage speaker diarization of broadcast news". In: *IEEE Transactions on Audio, Speech, and Language Processing* 14.5 (Sept. 2006), pp. 1505–1512. ISSN: 1558-7916. DOI: 10.1109/TASL.2006.878261.

# A    Code Listings

## A.1    main.py

```python
1  import tensorflow as tf
2
3  from scipy.io import wavfile
4  from textgrid import TextGrid
5  from diarization_methods.TfLearnNN import TfLearnNN
6  from os import walk
7  from os.path import splitext
8  from numpy import array, vstack, hstack, zeros, ceil
9  from numpy.random import shuffle, randint
10 import time
11
12 audio_sample_rate = 8000
13
14 # Parameters to be set in some more elegant way later on.
15 audio_files_path  = '../media/Full_Test_Files_8000/'
16 text_grids_path  = '../media/Full_Test_TextGrids/'
17 load_all_to_memory = True
18 window_size = int(audio_sample_rate * 0.2)
19 num_classes = 4
20 num_channels= 2
21 # training_flag  = True
22  training_flag  = False
23 epochs = 50
24 batch_size  = 100
25 percent_train  = .8
26 params = {
27     'num_channels' : 2,
28      'window_size' : window_size,
29      'num_classes' : 4,
30           'nfft' : 100
31 }
32 diarization_method = TfLearnNN(params)
33 load_path = '../media/tflearn_nn.save'
34
35 # run_audio_file_path  = '../media/Full Test Files 8000/HS_D01.wav'
36 #  text_grid_load_file_path  = '../media/Full Test TextGrids/HS_D01.TextGrid'
37 #  text_grids_save_file_path  = '../media/TexGrid Output/'
38
39 # save_path = '../media/tflearn_nn_2.save'
40 save_path = None
41 # load_path = None
42 load_path = '../media/tflearn_nn.save'
43
44 # End parameters section
45
46
47 def   get_file_list  (directory):
48     """
49     Used to get all  the  file  names in the directory  specified .
50     :param directory: the path to the  directory  to get  the  file  names from (ending in '/'
51     :return: a  list  of  all   files  names in directory  starting  with directory
52     """
```

```python
53      f = []
54      for (dirpath, dirnames, filenames) in walk(directory):
55          f = filenames
56          break
57      f.sort()
58      for i in range(len(f)):
59          f[i] = directory + f[i]
60      return f
61
62
63  def get_label(path):
64      """
65      Used to load a TextGrid from disk and return the array representation of it.
66      :pram path: Path to the TextGrid
67      :return: Array representation of the TextGrid.
68      """
69      grid = TextGrid(name=path)
70      grid.read(path, Fs=int(audio_sample_rate / window_size))
71      return grid.FsArrayCombined
72
73
74  class Trainer:
75      def __init__(self):
76          self.data_paths = None
77          self.label_paths = None
78          self.training_data = None
79          self.training_labels = None
80          self.testing_data = None
81          self.testing_labels = None
82          self.imax = -1
83
84          self.data_paths = array([get_file_list(audio_files_path)])
85          self.label_paths = array([get_file_list(text_grids_path)])
86          if len(self.data_paths) != len(self.label_paths):
87              print('Error! Data and labels do not match up!')
88              ValueError()
89
90          self.split_train_test()
91
92          if load_all_to_memory:
93              print('Loading Data...', end='', flush=True)
94              self.load_all_data_to_memory()
95              print('Done!')
96          else:
97              print("Training is going to take a bit longer without loading all data first ..")
98
99      def get_chunk(self, is_train=True):
100         data = self.training_data if is_train else self.testing_data
101         labels = self.training_labels if is_train else self.testing_labels
102
103         i = randint(0, len(data))
104         x = data[i]  # should be a list at this point
105         label = labels[i]
106
107         if not load_all_to_memory:  # means that data is a path instead
108             fs, x = wavfile.read(x[0])
109             label = get_label(label[0])
110         start = randint(0, len(label) - 1)
```

15

```python
111              return x.T[:, start * window_size:start * window_size + window_size], label[start]

112
113         def get_complete_file_chunk(self, is_train=True):
114             data = self.training_data if is_train else self.testing_data
115             labels = self.training_labels if is_train else self.testing_labels

116
117             i = randint(0, len(data))
118             x = data[i]   # should be a list at this point
119             label = labels[i]

120
121             if not load_all_to_memory:  # means that data is a path instead
122                 fs, x = wavfile.read(x[0])
123                 label = get_label(label[0])
124             start = randint(0, len(label) - 1)
125             return x.T[:, start * window_size:start * window_size + window_size], label[start]

126
127         def load_all_data_to_memory(self):
128             paths = array(self.training_data, copy=True)
129             self.training_data = []
130             for path in paths:
131                 fs, x = wavfile.read(path[0])
132                 self.training_data.append(x)

133
134             paths = array(self.testing_data, copy=True)
135             self.testing_data = []
136             for path in paths:
137                 fs, x = wavfile.read(path[0])
138                 self.testing_data.append(x)

139
140             paths = array(self.training_labels, copy=True)
141             self.training_labels = []
142             for path in paths:
143                 self.training_labels.append(get_label(path[0]))

144
145             paths = array(self.testing_labels, copy=True)
146             self.testing_labels = []
147             for path in paths:
148                 self.testing_labels.append(get_label(path[0]))

149
150         def split_train_test(self):
151             shuf = vstack((self.data_paths, self.label_paths)).T
152             shuffle(shuf)

153
154             self.imax = i_max = int(len(self.data_paths[0]) * percent_train) + 1
155             self.training_data = shuf[0:i_max, [0]]
156             self.training_labels = shuf[0:i_max, [1]]
157             self.testing_data = shuf[i_max - 1:-1, [0]]
158             self.testing_labels = shuf[i_max - 1:-1, [1]]
159             # At this point the training and testing sets contain a list of paths to data and labels.

160
161         def train(self):
162             try:
163                 if load_path is not None:
164                     print('Loading_Diarization')
165                     diarization_method.load(load_path)

166

167
168                 X = array([], dtype=float).reshape(0,window_size)
```

```python
169                Y = array([], dtype=int).reshape(0,1)
170                for i in range(epochs):
171                    for j in range(batch_size):
172                        x, label = self.get_chunk()
173                        X = vstack((X,x))
174                        Y = vstack((Y,label))
175                    diarization_method.train_on_data(X, Y)
176
177                print('Finished_Training...')
178                print('Testing...')
179
180                # Testing section
181                l = len(self.testing_data)
182                for i in range(l):
183                    x, label = self.get_chunk(is_train=False)
184                    X = vstack((X,x))
185                    Y = vstack((Y,label))
186                num_error = diarization_method.get_train_error(X, Y)
187                num_train = Y.shape[0]
188                training_accuracy = (num_train-num_error)/num_train*100
189                print("%d/%d_Mistakes._Training_Accuracy:_%.2f%%"%(int(num_error),num_train,
                          training_accuracy))
190
191            except Exception as ex:
192                print("ERROR!_And_exception_Occurred!")
193                print(ex)
194            except KeyboardInterrupt:
195                print("\n_Interrupted..")
196            finally:
197                if save_path is not None:
198                    print('Saving_Diarization_to_%s...'%(save_path))
199                    diarization_method.save(save_path)
200
201
202    class Predictor:
203        def __init__(self):
204            self.data_paths = run_audio_file_path
205            self.label_paths = None
206            self.data = []
207            self.labels = []
208            self.imax = -1
209            self.textgrid = None
210
211            if text_grids_path is not None:
212                self.labels = get_label(text_grid_load_file_path)[1:-1]
213
214            fs, x = wavfile.read(self.data_paths)
215            X0 = self.create_data_chunks(x[:,0], window_size,window_size)
216            X1 = self.create_data_chunks(x[:,1], window_size,window_size)
217            self.data = hstack((X0,X1))
218
219        def create_data_chunks(self,x,window_size,step_size):
220            window_size = int(window_size)
221            step_size = int(step_size)
222
223            # drop elements outside window
224            if int(len(x)%(step_size)) != 0:
225                x = x[:-int(len(x)%(step_size))]
```

```python
226
227          # ending
228          data_nb = int( ceil (( len(x) − window_size)/step_size))
229
230          start=0
231          end=window_size
232          data = zeros([data_nb,window_size])
233
234          for i in range(data_nb):
235              data[i ,:] = x[start :end]
236              start=start+step_size
237              end=end+step_size
238
239          return data
240
241      def run( self ):
242          try :
243              if load_path is not None:
244                  print('Loading_Diarization')
245                  diarization_method.load(load_path)
246
247              # diarization_method.run_on_data(self.data_paths)
248              num_error = diarization_method.get_train_error( self .data, self . labels )
249              num_train = self. labels .shape[0]
250              training_accuracy = (num_train−num_error)/num_train∗100
251              print("%d/%d_Mistakes._Training_Accuracy:_%.2f%%"%(int(num_error),num_train,
                     training_accuracy))
252
253          except Exception as ex:
254              print("ERROR!_and_exception_Occurred!")
255              print (ex)
256          except KeyboardInterrupt:
257              print("\n_Interrupted..")
258          finally :
259              print("Done")
260
261
262  if __name__ == '__main__':
263      # TODO: command line argument parsing. Use are parse
264      if training_flag is True:
265          trainer = Trainer()
266          trainer . train ()
267      else :
268          predictor = Predictor()
269          predictor .run()
```

## A.2   DiarizationBaseClass.py

```python
1  from abc import ABC, abstractmethod
2  import tensorflow as tf
3
4
5  class DiarizationBaseClass(ABC):
6      """
7      The DiarizationBaseClass provides the interfaces between custom diarization
8      methods and the wrapper code. It provides a tensorflow session available as
9      ' self . sess ' that is usable within the 'train_on_data ', ' get_train_error ', and
```

```python
10        'run_on_data' methods.
11
12        You may add other variables needed within your own implementation. Not here
13        please. You can pass parameters into your custom diarization method via the
14        'params' init parameter. This can be an array, dictonary, or object. In
15        other words, that is the only parameter you'll ever need. You can pack
16        everything into it.
17        """
18
19    def __init__ (self, params=None):
20        """
21        Initialization for the base diarization class.
22        :param params: the tuple, dictonary, or object used to initialize the
23                        diarization method.
24        """
25        self. tf_inititializer   = tf. global_variables_initializer ()
26        self.sess = None
27        self.init_diarization_method(params)
28        pass
29
30    @abstractmethod
31    def init_diarization_method(self, params):
32        """
33        This method is called once to initialize the speech diarization algorithm.
34        :param params: the tuple, dictonary, or object used to initialize the
35                        diarization method.
36        """
37        pass
38
39    @abstractmethod
40    def train_on_data(self, data, label):
41        """
42        This method is meant to be called once per iteration during the training cycle.
43
44        The tensorflow session is available as 'self.sess'. tensorflow session initialization
45        has already been taken care of with the 'tf. global_variables_initializer ()'
46
47        :param data: This is the single set of data to be used by a single training
48                     operation. This will change with each call of 'train_on_data'.
49        :param label: This is the textgrid label for this data sample.
50        """
51        pass
52
53    @abstractmethod
54    def run_on_data(self, data):
55        """
56        This method will be used to run the diarization method on a set of data and
57        returns the TextGrid result.
58        :param data: the audiofile to run the diarization on.
59        :return: the TextGrid object.
60        """
61        pass
62
63    @abstractmethod
64    def get_train_error (self, test_data, test_label):
65        """
66        This method will be called at the end of each epoch to get the error of the
67        diarization method.
```

```
68        :param test_data: the data to use to evaluate the accuracy of the model.
69        :param test_label: the label to use to evaluate the accuracy of the model.
70        :return: The error of the diarization method.
71        """
72
73    @abstractmethod
74    def load( self , path):
75        """
76        This method will be called when we want to load the diarization method from disk
77        :param path: the path to the storage file that is to be loaded.
78        :return: True on success, False on failure
79        """
80        pass
81
82    @abstractmethod
83    def save( self , path):
84        """
85        This method will be called when we want to save the diarization method to disk to
86        be recalled later .
87        :param path: the path to where the storage file should be put.
88        :return: True on success, False on failure
89        """
```

## A.3   TFLearnNN.py

```
1  from diarization_methods.DiarizationBaseClass import DiarizationBaseClass
2  import tensorflow as tf
3  import tflearn
4  from tflearn . layers .core import input_data, dropout, fully_connected
5  from tflearn . layers .conv import conv_2d, max_pool_2d
6  from tflearn . layers .normalization import local_response_normalization
7  from tflearn . layers .estimator import regression
8  import time
9  import numpy as np
10
11  from numpy import array, zeros
12
13
14  class  TfLearnNN(DiarizationBaseClass):
15      def run_on_data(self, data):
16          # create fft
17          if len(test_data) != len( test_label ):
18              test_data = test_data.reshape(len( test_label ),−1)
19
20          Xfft = self . __get_fft (test_data)
21          X = np.hstack((test_data,Xfft))
22
23          # create one−hot label
24          Y = self. label_to_one_hot ( test_label , self .num_classes)
25
26          yhat = self .model.predict(X)
27          yhat = np.argmax(yhat, axis=1)
28
29          return  self . __validate_results (yhat, test_label )
30
31      def label_to_one_hot ( self , label , num_classes):
32          """
```

```python
33              converts
34              into
35              :pram label:    [0,  3,  2,  1,   ...]
36              :return:        [[1,  0,  0,  0,  ...],
37                               [0,  0,  0,  1,  ...],
38                               [0,  0,  1,  0,  ...],
39                               [0,  1,  0,  0,  ...]]
40              """
41              y = np.eye(self.num_classes)[label.astype(int)].reshape(-1,num_classes)
42              return y
43
44          def init_diarization_method(self, params):
45              self.window_size  = params['window_size']
46              self.nfft         = params['nfft']
47              self.num_channels = params['num_channels']
48              self.num_classes  = params['num_classes']
49              self.total_width  = (self.window_size + self.nfft)*self.num_channels
50
51              input_layer = tflearn.input_data(shape=[None, self.total_width], name='input')
52              dense1 = tflearn.fully_connected(input_layer, 300, activation='relu')
53
54              dense2 = tflearn.fully_connected(dense1, 200, activation='relu')
55              dropout2 = tflearn.dropout(dense2, 0.8)
56
57              dense3 = tflearn.fully_connected(dropout2, 200, activation='relu')
58              dropout3 = tflearn.dropout(dense3, 0.8)
59
60              dense4 = tflearn.fully_connected(dropout3, 500, activation='relu')
61              dropout4 = tflearn.dropout(dense4, 0.8)
62
63              softmax = tflearn.fully_connected(dropout4, 4, activation='softmax')
64
65              self.net = tflearn.regression(softmax, loss='categorical_crossentropy',name='target')
66              self.model = tflearn.DNN(self.net, tensorboard_verbose=2)
67
68          def train_on_data(self, data, label):
69              if len(data) != len(label):
70                  data = data.reshape(len(label),-1)
71                  test_data = test_data.reshape(len(test_label),-1)
72
73              # create fft
74              Xfft = self.__get_fft(data)
75              X = np.hstack((data,Xfft))
76
77              # create one-hot label
78              Y = self.label_to_one_hot(label, self.num_classes)
79
80              self.model.fit(X, Y, n_epoch=1, show_metric=True, \
81                             validation_set =(X, Y),
82                             snapshot_epoch=False,run_id='TfLearnNNk')
83
84          def get_train_error(self, test_data, test_label):
85              # create fft
86              if len(test_data) != len(test_label):
87                  test_data = test_data.reshape(len(test_label),-1)
88
89              Xfft = self.__get_fft(test_data)
90              X = np.hstack((test_data,Xfft))
```

```
91
92          # create one−hot label
93          Y = self. label_to_one_hot ( test_label , self .num_classes)
94
95          yhat = self .model.predict(X)
96          yhat = np.argmax(yhat, axis=1)
97
98          return  self .  __validate_results (yhat, test_label )
99
100     def load( self ,  path):
101          self .model.load(path)
102
103     def save( self ,  path):
104          self .model.save(path)
105
106     def  __validate_results ( self ,  yhat, y):
107          num_error = 0
108          for  i  in  range(len(yhat)):
109              if  yhat[i ]  != y[i ]:
110                  num_error += 1
111          return num_error
112
113     def  __get_fft ( self ,  data):
114          if  self .num_channels is 2:
115              Xfft0 = abs(np.fft. fft (data [:,: self .window_size], self . nfft ))
116              Xfft1 = abs(np.fft. fft (data [:, self .window_size:], self . nfft ))
117              Xfft = np.hstack((Xfft0,Xfft1))
118          else :
119              Xfft = abs(np.fft. fft (data [:, self .window_size:], self . nfft ))
120          return  Xfft
```

## A.4   plot.py

```
1 import matplotlib.pyplot as  plt
2 import numpy as np
3 from scipy. io  import  wavfile
4 import scipy. fftpack  as  fftpack
5 from scipy  import  signal
6 import sounddevice as sd
7 import matplotlib.animation as animation
8 from matplotlib  import  rcParams
9
10 rcParams.update({'figure.autolayout': True})
11
12
13 def track_audio( fs ,  start=0, end="end"):
14     Y_MIN = −2
15     Y_MAX = 2
16
17     x = np.arange(start, end + 1, 0.01);
18
19     def update_line(num, line):
20         i = x[num]
21         line . set_data ([ i ,  i ],  [Y_MIN, Y_MAX])
22         return  line ,
23
24     l,  _ = plt.plot( start ,  −1, end, 1,  linewidth=2, color='red')
```

```python
25      line_anim = animation.FuncAnimation(fig, update_line, len(x), fargs=(l,), interval=1 / fs, blit
            =True, repeat=False)
26      plt.show()
27

28
29  def plot_spectrogram(x, fs, start=0, end="end"):
30      if end is "end":
31          Pxx, freqs, bins, im = plt.specgram(x[start * fs:-1],
32                                          NFFT=512, Fs=fs, noverlap=100, cmap=plt.cm.gist_heat
                                                )
33          end = max(bins)
34
35      else:
36          Pxx, freqs, bins, im = plt.specgram(x[start * fs:start * fs + end * fs],
37                                          NFFT=512, Fs=fs, noverlap=100, cmap=plt.cm.gist_heat
                                                )
38          if end > max(bins):
39              end = max(bins)
40
41      plt.ylim(0, max(freqs))
42      plt.xlim([start, end])
43      plt.ylabel("Frequency_(Hz)")
44      plt.xlabel("Time_(Sec)")
45      # plt.colorbar()
46

47
48  def play_audio(data, fs, start=0, end="end", blocking=True):
49      if end is "end":
50          sd.play(data[int(start * fs):-1], fs, blocking=blocking)
51      else:
52          sd.play(data[start * fs:int(start * fs + end * fs)], fs, blocking=blocking)
53

54
55  def plot_audio(x, fs, start=0, end="end"):
56      if end is "end":
57          xaxis = np.linspace(start, len(x) / fs, num=len(x[start * fs:]))
58          plt.plot(xaxis, x[start * fs:] / float(max(x)), linewidth=0.25)
59          end = len(x) / float(fs)
60
61      else:
62          if end > len(x) / fs:
63              end = len(x) / fs
64          xaxis = np.linspace(start, end, num=len(x[int(start * fs):int(end * fs)]))
65          plt.plot(xaxis, x[int(start * fs):int(end * fs)] / float(max(x)), color='#3030e0',
                linewidth=0.15)
66
67      plt.ylim(-1.2, 1.2)
68      plt.xlim([start, end])
69      plt.ylabel("Magnitude")
70      plt.xlabel("Time_(Sec)")
71

72
73  def plot_bounds_lines(changes, marks, start=0, end="end"):
74      x = changes
75      if end is "end":
76          last = x.shape[0] + 1
77      else:
78          last = np.argmax(x > end)
```

23

```python
79
80      data_range = range(np.argmin(x < start), last)
81
82      if max(marks) > 2:
83          plt.plot((-1, -1), (-1, -1), 'b-', linewidth=2, label="Speaker_1")
84          plt.plot((-1, -1), (-1, -1), 'r-', linewidth=2, label="Speaker_2")
85          plt.plot((-1, -1), (-1, -1), 'g-', linewidth=2, label="Both")
86      else:
87          plt.plot((-1, -1), (-1, -1), 'b-', linewidth=2, label="Speech")
88
89      plt.legend(loc=2, fancybox=True, framealpha=0.8, prop={'size': 9})
90
91      for j, mark in enumerate(marks):
92          if mark == 1:
93              plt.plot((x[j], x[j + 1]), (-1.05, -1.05), 'b-', linewidth=2)
94              plt.plot((x[j + 1], x[j + 1]), (-1.03, -1.07), 'k-', linewidth=2, zorder=10)
95
96          elif mark == 2:
97              plt.plot((x[j], x[j + 1]), (-1.05, -1.05), 'r-', linewidth=2)
98              plt.plot((x[j + 1], x[j + 1]), (-1.03, -1.07), 'k-', linewidth=2, zorder=10)
99
100         elif mark == 3:
101             plt.plot((x[j], x[j + 1]), (-1.05, -1.05), 'g-', linewidth=2)
102             plt.plot((x[j + 1], x[j + 1]), (-1.03, -1.07), 'k-', linewidth=2, zorder=10)
103
104         else:
105             pass
106
107     plt.ylim(-1.2, 1.2)
108     plt.xlim([start, end])
109     plt.xlabel("Time_(Sec)")
110
111
112 def plot_bounds_fill(changes, marks, start=0, end="end"):
113     x = changes
114     if end is "end":
115         last = x.shape[0] + 1
116     else:
117         last = np.argmax(x >= end) + 1
118
119     data_range = range(np.argmin(x < start), last)
120
121     # create legend information
122     # plt.axvspan(0, 0, alpha=0.5, color='b', label="\"N\" Human")
123     # plt.axvspan(0, 0, alpha=0.5, color='r', label="\"S\" Human")
124     # plt.axvspan(0, 0, alpha=0.5, color='g', label="\"S\" Human")
125     plt.legend(loc=2, fancybox=True, framealpha=0.8, prop={'size': 9})
126
127     for j, mark in enumerate(marks):
128         if mark == 1:
129             plt.axvspan(x[j], x[j + 1], alpha=0.2, color='b')
130
131         elif mark == 2:
132             plt.axvspan(x[j], x[j + 1], alpha=0.2, color='r')
133
134         elif mark == 3:
135             plt.axvspan(x[j], x[j + 1], alpha=0.2, color='g')
136
```

```
137            else :
138                pass
139
140
141  def  plot_fft () :
142      N = 600 # sample points
143      T = 1 / 800.0  # sample spacing
144      x = np.linspace(0,  N ∗ T,  N)
145      y = np.sin(50.0 ∗ 2.0 ∗ np.pi ∗ x) + 0.5 ∗ np.sin(80.0 ∗ 2.0 ∗ np.pi ∗ x)
146      yf = fftpack. fft (y)
147      xf = np.linspace(0.0,  1.0 / (2.0 ∗ T), N / 2)
148      fig ,  ax = plt.subplots()
149      ax.plot(xf,  2.0 / N ∗ np.abs(yf[:N // 2]))
150
151
152  if  __name__ == '__main__':
153      main()
```

## A.5   textgrid.py

An existing TextGrid parser was used from `https://github.com/kylebgorman/`
`textgrid/blob/master/textgrid/textgrid.py`. Three additional functions were added
to assist in neural network training:

```
1  def get_bin_changes(x, Fs):
2      FsTimeChanges = [0]
3      FsChangeMarks = [x[0]]
4      for i in range(len(x)):
5          if x[i] != x[i − 1]:
6              FsTimeChanges.append(i)
7              FsChangeMarks.append(x[i])
8
9      FsTimeChanges.append(i)
10     FsTimeChanges = np.array(FsTimeChanges)∗1/float(Fs)
11     FsChangeMarks = np.array(FsChangeMarks)
12     return FsTimeChanges, FsChangeMarks
```

```
1  class  IntervalTier (object):
2  ...
3    def  tier_to_array ( self ,  Fs,  maxTime):
4      array = np.zeros(int(np. ceil (Fs ∗ maxTime)))
5      current_interval  = 0
6      for i  in range(array.shape[0]) :
7        if  i/float (Fs) < self [ current_interval ]. maxTime:
8          if  self [ current_interval ]. mark == 'N':
9            array[i] = 1
10         else :
11       current_interval  += 1
12
13     self .FsArray = array
14     self .FsTimeChanges, self.FsChangeMarks = get_bin_changes(array, Fs)
15  ...
```

```
1  class  TextGrid(object):
2  ...
3
4    def  combine_grids( self ):
```

```python
5         self .FsArrayCombined = self.tiers[0].FsArray + self. tiers [1]. FsArray * 2
6         self .FsTimeChangesCombined, self.FsChangeMarksCombined = get_bin_changes(self.
              FsArrayCombined, self.Fs)
7  ...
```

```python
1  class  TextGrid(object):
2  ...
3    def read( self , f, round_digits=DEFAULT_TEXTGRID_PRECISION, Fs=None):
4    """"
5    Read the tiers  contained in  the Praat−formatted TextGrid file
6    indicated  by string  f. Times are rounded to the specified  precision .
7    """"
8     self .Fs = Fs
9    encoding = detectEncoding(f)
10   with codecs.open(f, 'r', encoding=encoding) as source:
11     source. readline ()   # header junk
12     source. readline ()   # header junk
13     source. readline ()   # header junk
14
15     self .minTime = round(float(source.readline().split () [2]) ,  round_digits)
16     self .maxTime = round(float(source.readline().split () [2]) ,  round_digits)
17     source. readline ()   # more header junk
18     m = int(source.readline(). rstrip (). split () [2])   # will be self .n
19     source. readline ()
20     for  i  in  range(m):  # loop over grids
21       source. readline ()
22       if source.readline (). rstrip (). split () [2]  == '"IntervalTier"':
23         inam = source.readline(). rstrip (). split ('␣=␣') [1]. strip ('"')
24         imin = round(float(source.readline (). rstrip (). split () [2]) ,  round_digits)
25         imax = round(float(source.readline (). rstrip (). split () [2]) ,  round_digits)
26         itie  = IntervalTier(inam)
27         for  j  in  range(int(source. readline (). rstrip (). split () [3]) ):
28           source.readline (). rstrip (). split ()   # header junk
29           jmin = round(float(source.readline (). rstrip (). split () [2]) ,  round_digits)
30           jmax = round(float(source.readline(). rstrip (). split () [2]) ,  round_digits)
31           jmrk = _getMark(source)
32           if  jmin < jmax: # non−null
33               itie .addInterval( Interval (jmin, jmax, jmrk))
34         self .append(itie)
35       else :   # pointTier
36         inam = source.readline(). rstrip (). split ('␣=␣') [1]. strip ('"')
37         imin = round(float(source.readline (). rstrip (). split () [2]) ,  round_digits)
38         imax = round(float(source.readline(). rstrip (). split () [2]) ,  round_digits)
39         itie  = PointTier(inam)
40         n = int(source. readline (). rstrip (). split () [3])
41         for  j  in  range(n):
42           source. readline (). rstrip ()   # header junk
43           jtim  = round(float(source.readline (). rstrip (). split () [2]) ,round_digits)
44           jmrk = _getMark(source)
45            itie .addPoint(Point(jtim, jmrk))
46         self .append(itie)
47
48    if  Fs != None:
49      for  tier  in  self . tiers :
50        tier . tier_to_array (Fs,  self .maxTime)
51      if  len( self . tiers ) == 2:
52        self .combine_grids()
53  ...
```