

# 4 Deep Learning Models for Speech Recognition

## Data Science - Capstone Project Submission

- Student Name: **James Toop**
- Student Pace: **Self Paced**
- Scheduled project review date/time: **29th October 2021 @ 21:30 BST**
- Instructor name: **Jeff Herman / James Irving**
- Blog URL: <https://toopster.github.io/> (<https://toopster.github.io/>)

In [1]:

```
1 # Import relevant libraries and modules for creating and training networks
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import wave
7 import soundfile as sf
8 import librosa, librosa.display
9 import IPython.display as ipd
10 import os
11 import json
12
13 import tensorflow as tf
14 from tensorflow.keras.layers.experimental import preprocessing
15 from tensorflow.keras import layers
16 from tensorflow.keras import models
17 from sklearn.model_selection import train_test_split
18 from sklearn.preprocessing import LabelEncoder
19
20 import logging
21 logging.getLogger("tensorflow").setLevel(logging.ERROR)
22
23 import pathlib
24 from pathlib import Path
25
26 import shared_functions.preprocessing as preprocess
```

In [2]:

```
1 # Set seed for reproducibility
2 seed = 123
3 tf.random.set_seed(seed)
4 np.random.seed(seed)
```

In [3]:

```
1 def visualise_results(results):
2     '''
3     Visualise the results from the model history plotting training and
4     validation accuracy and loss vs. epoch
5
6     Params:
7         results: Model history
8     '''
9     history = results.history
10
11     plt.figure(figsize=(20,8))
12     plt.xticks(fontsize=12)
13     plt.yticks(fontsize=12)
14
15     plt.subplot(1, 2, 1)
16     plt.plot(history['val_loss'])
17     plt.plot(history['loss'])
18     plt.legend(['Validation Loss', 'Training Loss'], fontsize=12)
19     plt.title('Loss', fontsize=18)
20     plt.xlabel('Epochs', fontsize=14)
21     plt.ylabel('Loss', fontsize=14)
22
23     plt.subplot(1, 2, 2)
24     plt.plot(history['val_acc'])
25     plt.plot(history['acc'])
26     plt.legend(['Validation Accuracy', 'Training Accuracy'], fontsize=12)
27     plt.title('Accuracy', fontsize=18)
28     plt.xlabel('Epochs', fontsize=14)
29     plt.ylabel('Accuracy', fontsize=14)
30     plt.show()
```

In [4]:

```
1 # Create the train, test and validation datasets for the Speech Commands dataset
2 sc_data_path = "speech_commands_data.json"
3 (
4     sc_X_train,
5     sc_y_train,
6     sc_X_val,
7     sc_y_val,
8     sc_X_test,
9     sc_y_test,
10 ) = preprocess.create_train_test(sc_data_path, "MFCCs")
```

Datasets loaded...

In [5]:

```
1 # Create the train, test and validation datasets for the Ultrasuite Top 35 data
2 us_data_path = 'ultrasuite_top35_data.json'
3 (
4     us_X_train,
5     us_y_train,
6     us_X_val,
7     us_y_val,
8     us_X_test,
9     us_y_test,
10 ) = preprocess.create_train_test(us_data_path, 'MFCCs')
```

Datasets loaded...

In [6]:

```
1 def check_array_shapes(X_train, y_train, X_val, y_val, X_test, y_test):
2     '''
3     Output the training, test and validation dataset shapes
4
5     Params:
6         X_train (ndarray): Inputs for the training dataset
7         y_train (ndarray): Targets for the training dataset
8         X_val (ndarray): Inputs for the validation dataset
9         y_val (ndarray): Targets for the validation dataset
10        X_test (ndarray): Inputs for the test dataset
11        y_test (ndarray): Targets for the test dataset
12    '''
13    print ("Number of training samples: " + str(X_train.shape[0]))
14    print ("Number of testing samples: " + str(X_test.shape[0]))
15    print ("Number of validation samples: " + str(X_val.shape[0]))
16    print ("X_train shape: " + str(X_train.shape))
17    print ("y_train shape: " + str(y_train.shape))
18    print ("X_test shape: " + str(X_test.shape))
19    print ("y_test shape: " + str(y_test.shape))
20    print ("X_val shape: " + str(X_val.shape))
21    print ("y_val shape: " + str(y_val.shape))
```

In [7]:

```
1 def reformat_y(y):
2     '''
3     Reformats / One Hot Encodes targets
4
5     Params:
6         y (ndarray): Input targets
7
8     Returns:
9         y (ndarray): One hot encoded targets
10    '''
11    y = LabelEncoder().fit_transform(y)
12    y = tf.keras.utils.to_categorical(y)
13    return y
```

In [8]:

```
1 # Check the array shapes for the Ultrasuite dataset
2 check_array_shapes(us_X_train,
3                     us_y_train,
4                     us_X_val,
5                     us_y_val,
6                     us_X_test,
7                     us_y_test)
```

Number of training samples: 3942  
Number of testing samples: 1233  
Number of validation samples: 986  
X\_train shape: (3942, 44, 13, 1)  
y\_train shape: (3942,)  
X\_test shape: (1233, 44, 13, 1)  
y\_test shape: (1233,)  
X\_val shape: (986, 44, 13, 1)  
y\_val shape: (986,)

In [9]:

```
1 us_y_test[:10]
```

Out[9]:

```
array([23, 31, 19,  6, 22, 10, 14,  3,  6, 11])
```

In [10]:

```
1 # One-hot encode Ultrasuite Top 35 labels
2 us_train_y = reformat_y(us_y_train)
3 us_test_y = reformat_y(us_y_test)
4 us_val_y = reformat_y(us_y_val)
```

```
1 us_test_y[:10]
```

## 4.1 Model 1: Create a simple baseline model

In [12]:

```
1 def build_baseline_model(input_shape,
2                           output_units,
3                           loss_func='categorical_crossentropy',
4                           learning_rate=0.0001):
5     '''
6     Build a baseline model
7
8     Params:
9         input_shape (tuple): Shape of array representing a sample train
10        output_units (int): Number of targets / categories
11        loss_func (str): Loss function to use
12        learning_rate (float): Learning rate
13
14    Returns:
15        baseline_model: Tensorflow model
16    '''
17    baseline_model = tf.keras.models.Sequential()
18    baseline_model.add(tf.keras.layers.InputLayer(input_shape=input_shape))
19    baseline_model.add(tf.keras.layers.Flatten())
20    baseline_model.add(tf.keras.layers.BatchNormalization())
21    baseline_model.add(tf.keras.layers.Dense(output_units, activation='softmax'))
22
23    # Set optimizer and learning rate
24    optimiser = tf.optimizers.Adam(learning_rate=learning_rate)
25
26    # Compile the baseline model
27    baseline_model.compile(loss=loss_func,
28                          optimizer=optimiser,
29                          metrics=['acc'])
30
31    # Print summary for model
32    baseline_model.summary()
33
34    return baseline_model
```

In [13]:

```
1 # Function for fitting the model
2 def fit_model(model,
3               epochs,
4               batch_size,
5               patience,
6               X_train,
7               y_train,
8               X_val,
9               y_val):
10     '''
11     Fit the model
12
13     Params:
14         model : Input Tensorflow model
15         epochs (int): Number of training epochs
16         batch_size (int): Number of samples per batch
17         patience (int): Number of epochs to wait before early stop,
18                        if there no improvement on accuracy
19         X_train (ndarray): Inputs for the training dataset
20         y_train (ndarray): Targets for the training dataset
21         X_val (ndarray): Inputs for the validation dataset
22         y_val (ndarray): Targets for the training dataset
23
24     Returns:
25         results: Training history
26     '''
27     # Define early stopping criteria
28     early_stopping = tf.keras.callbacks.EarlyStopping(monitor='accuracy',
29                                                       min_delta=0.001,
30                                                       patience=patience)
31
32     # Fit the model
33     results = model.fit(X_train,
34                        y_train,
35                        epochs=epochs,
36                        batch_size=batch_size,
37                        validation_data=(X_val, y_val),
38                        callbacks=[early_stopping])
39     return results
```

In [14]:

```
1 # Function to save the model if so required
2 def save_model(save_model, save_path):
3     '''
4     Save the model
5
6     Params:
7         save_model : Input Tensorflow model
8         save_path (str): Path to save model including file extension .h5
9     '''
10     save_model.save(save_path)
```

#### 4.1.1 Baseline model for the Speech Commands dataset

In [15]:

```
1 # One-hot encode Speech Commands labels
2 sc_train_y = reformat_y(sc_y_train)
3 sc_test_y = reformat_y(sc_y_test)
4 sc_val_y = reformat_y(sc_y_val)
5
6 # Create baseline model for Speech Commands dataset
7 sc_input_shape = (sc_X_train[0].shape)
8 sc_output_units = 35
9 sc_baseline_model = build_baseline_model(sc_input_shape,
10                                         sc_output_units,
11                                         learning_rate=0.0001)
12
13 # Fit model
14 sc_epochs = 40
15 sc_batch_size = 32
16 sc_patience = 5
17 sc_baseline_results = fit_model(sc_baseline_model,
18                                 sc_epochs,
19                                 sc_batch_size,
20                                 sc_patience,
21                                 sc_X_train,
22                                 sc_train_y,
23                                 sc_X_val,
24                                 sc_val_y)
```

n it as-is.

Please report this to the TensorFlow team. When filing the bug, set the verbosity to 10 (on Linux, `export AUTOGRAPH\_VERBOSITY=10`) and attach the full output.

Cause: 'arguments' object has no attribute 'posonlyargs'

To silence this warning, decorate the function with @tf.autograph.exp

erimental.do\_not\_convert

1908/1908 [=====] - 8s 4ms/step - loss: 3.0156 - acc: 0.2003 - val\_loss: 2.4874 - val\_acc: 0.3122

Epoch 2/40

1908/1908 [=====] - 6s 3ms/step - loss: 2.3892 - acc: 0.3389 - val\_loss: 2.2717 - val\_acc: 0.3737

Epoch 3/40

1908/1908 [=====] - 6s 3ms/step - loss: 2.2549 - acc: 0.3775 - val\_loss: 2.1956 - val\_acc: 0.3955

Epoch 4/40

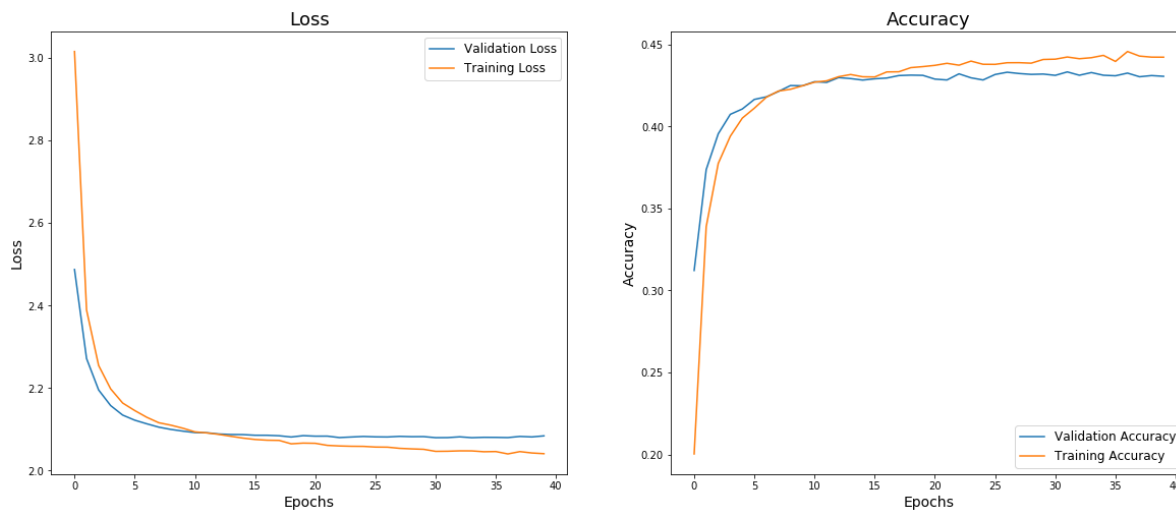
1908/1908 [=====] - 6s 3ms/step - loss: 2.1982 - acc: 0.3940 - val\_loss: 2.1577 - val\_acc: 0.4074

Epoch 5/40



In [16]:

```
1 # Visualise the training and validation loss and accuracy across epochs
2 visualise_results(sc_baseline_results)
```



In [17]:

```
1 # Evaluate the training results
2 sc_baseline_train = sc_baseline_model.evaluate(sc_X_train, sc_train_y)
3 sc_baseline_train
```

```
1908/1908 [=====] - 4s 2ms/step - loss: 1.947
3 - acc: 0.4745
```

Out[17]:

```
[1.947283148765564, 0.47454628348350525]
```

In [18]:

```
1 # Evaluate the test results
2 sc_baseline_test = sc_baseline_model.evaluate(sc_X_test, sc_test_y)
3 sc_baseline_test
```

```
597/597 [=====] - 1s 2ms/step - loss: 2.0892
- acc: 0.4292A: 0s - loss: 2.0889 - acc: 0.4
```

Out[18]:

```
[2.0892109870910645, 0.4291629493236542]
```

## 4.1.2 Baseline model for the Ultrasuite dataset

In [19]:

```
1 # Create baseline model for Ultrasuite dataset
2 us_input_shape = (us_X_train[0].shape)
3 us_output_units = 35
4 us_baseline_model = build_baseline_model(us_input_shape,
5                                         us_output_units,
6                                         learning_rate=0.0001)
7
8 # Fit model
9 us_epochs = 40
10 us_batch_size = 32
11 us_patience = 5
12 us_baseline_results = fit_model(us_baseline_model,
13                                 us_epochs,
14                                 us_batch_size,
15                                 us_patience,
16                                 us_X_train,
17                                 us_train_y,
18                                 us_X_val,
19                                 us_val_y)
```

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
=====		
flatten_1 (Flatten)	(None, 572)	0
=====		
batch_normalization_1 (Batch Normalization)	(None, 572)	2288
=====		
dense_1 (Dense)	(None, 35)	20055
=====		

Total params: 22,343

Trainable params: 21,199

Non-trainable params: 1,144

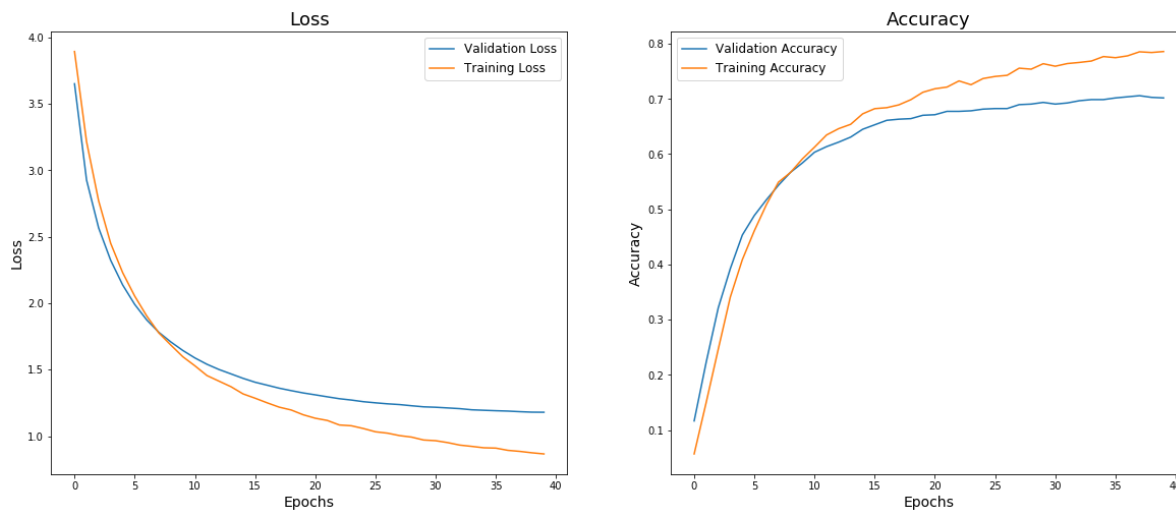
Epoch 1/40

WARNING: AutoGraph could not transform <function Model.make\_train\_function.<locals>.train\_function at 0x7face67fd8c0> and will run it as-is.

Please report this to the TensorFlow team. When filing the bug, set t

In [20]:

```
1 # Visualise the training and validation loss and accuracy across epochs
2 visualise_results(us_baseline_results)
```



In [21]:

```
1 # Evaluate the training results
2 us_baseline_results_train = us_baseline_model.evaluate(us_X_train, us_train_y)
3 us_baseline_results_train
```

```
124/124 [=====] - 0s 3ms/step - loss: 0.7894
- acc: 0.8130
```

Out[21]:

```
[0.7894143462181091, 0.8130390644073486]
```

In [22]:

```
1 # Evaluate the test results
2 us_baseline_results_test = us_baseline_model.evaluate(us_X_test, us_test_y)
3 us_baseline_results_test
```

```
39/39 [=====] - 0s 3ms/step - loss: 1.2542 -
acc: 0.6959
```

Out[22]:

```
[1.254226803779602, 0.6958637237548828]
```

### 4.1.3 Model Conclusion

It's starting point with the training accuracy for the Ultrasuite dataset being 81%. The model is clearly overfitting the training data and is not generalising well when shown unseen data given the validation accuracy of just over 67%.

What is of particular note is the vastly different accuracies when compared to the Speech Commands dataset which are both under 50%.

This could suggest that a more simple model works better for audio samples from children with a speech sound disorder.

## 4.2 Model 2: Baseline model with increased learning rate and batch size

In [23]:

```
1  # Create second baseline model on Ultrasuite dataset changing the learning rate
2  us_input_shape = (us_X_train[0].shape)
3  us_output_units = 35
4  us_baseline_model_2 = build_baseline_model(us_input_shape,
5                                             us_output_units,
6                                             learning_rate=0.001)
7
8
9  # Fit model
10 us_epochs = 40
11 us_batch_size = 64
12 us_patience = 5
13 us_baseline_results_2 = fit_model(us_baseline_model_2,
14                                   us_epochs,
15                                   us_batch_size,
16                                   us_patience,
17                                   us_X_train,
18                                   us_train_y,
19                                   us_X_val,
20                                   us_val_y)
```

```
=====
Total params: 22,343
Trainable params: 21,199
Non-trainable params: 1,144
```

Epoch 1/40

WARNING: AutoGraph could not transform <function Model.make\_train\_function.<locals>.train\_function at 0x7fadeaceb290> and will run it as-is.

Please report this to the TensorFlow team. When filing the bug, set the verbosity to 10 (on Linux, `export AUTOGRAPH\_VERBOSITY=10`) and attach the full output.

Cause: 'arguments' object has no attribute 'posonlyargs'

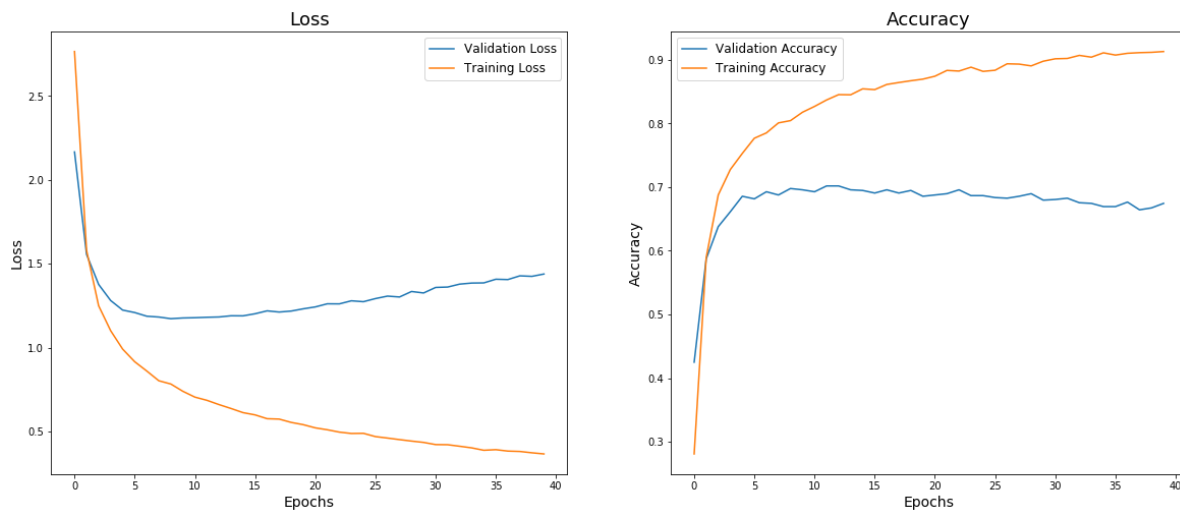
To silence this warning, decorate the function with @tf.autograph.experimental.do\_not\_convert

51/62 [=====>.....] - ETA: 0s - loss: 2.9259 - acc: 0.2377WARNING: AutoGraph could not transform <function Model.make\_test\_function.<locals>.test\_function at 0x7face68ede60> and will run it as-is.

Please report this to the TensorFlow team. When filing the bug, set t

In [24]:

```
1 # Visualise the training and validation loss and accuracy across epochs
2 visualise_results(us_baseline_results_2)
```



In [25]:

```
1 # Evaluate the training results
2 us_baseline_results_2_train = us_baseline_model_2.evaluate(us_X_train, us_train_y)
3 us_baseline_results_2_train
```

```
124/124 [=====] - 0s 2ms/step - loss: 0.2870
- acc: 0.9439
```

Out[25]:

```
[0.28699809312820435, 0.9439370632171631]
```

In [26]:

```
1 # Evaluate the test results
2 results_3_test = us_baseline_model_2.evaluate(us_X_test, us_test_y)
3 results_3_test
```

```
39/39 [=====] - 0s 4ms/step - loss: 1.5568 -
acc: 0.6545
```

Out[26]:

```
[1.5567854642868042, 0.65450119972229]
```

## 4.2.1 Model Conclusion

We have improved the training accuracy of the baseline model, to over 94%, by increasing the learning rate and increasing the batch size but this has only served to exacerbate the issue of overfitting with test accuracy not reducing to 65% and the difference between the two plots being dramatic.

In an effort to deal with the issue of overfitting, we will introduce some regularization layers and see what effect they have on performance.

## 4.3 Model 3: Adding hidden layers to the baseline model, deepening the network

In [27]:

```
1  # Build a second baseline model
2  def build_model_2(input_shape,
3                    output_units,
4                    loss_func='categorical_crossentropy',
5                    learning_rate=0.0001):
6
7      model_2 = tf.keras.models.Sequential()
8      model_2.add(tf.keras.layers.InputLayer(input_shape=input_shape))
9      model_2.add(tf.keras.layers.Flatten())
10     model_2.add(tf.keras.layers.BatchNormalization())
11     model_2.add(tf.keras.layers.Dense(256, activation='relu'))
12     # model_2.add(tf.keras.layers.Dropout(0.3))
13     model_2.add(tf.keras.layers.Dense(128, activation='relu'))
14     model_2.add(tf.keras.layers.Dense(64, activation='relu'))
15     model_2.add(tf.keras.layers.Dense(output_units, activation='softmax'))
16
17     # Set optimizer and learning rate
18     optimiser = tf.optimizers.Adam(learning_rate=learning_rate)
19
20     # Compile the baseline model
21     model_2.compile(loss=loss_func,
22                   optimizer=optimiser,
23                   metrics=['acc'])
24
25     # Print summary for model
26     model_2.summary()
27
28     return model_2
```

In [28]:

```
1 # Create baseline model for Ultrasuite dataset
2 us_input_shape = (us_X_train[0].shape)
3 us_output_units = 35
4 us_model_2 = build_model_2(us_input_shape,
5                             us_output_units,
6                             learning_rate=0.0001)
7
8 # Fit model
9 us_epochs = 40
10 us_batch_size = 64
11 us_patience = 3
12 us_results_3 = fit_model(us_model_2,
13                           us_epochs,
14                           us_batch_size,
15                           us_patience,
16                           us_X_train,
17                           us_train_y,
18                           us_X_val,
19                           us_val_y)
```

Model: "sequential\_3"

Layer (type)	Output Shape	Param #
=====		
flatten_3 (Flatten)	(None, 572)	0
<hr/>		
batch_normalization_3 (Batch Normalization)	(None, 572)	2288
<hr/>		
dense_3 (Dense)	(None, 256)	146688
<hr/>		
dense_4 (Dense)	(None, 128)	32896
<hr/>		
dense_5 (Dense)	(None, 64)	8256
<hr/>		
dense_6 (Dense)	(None, 35)	2275
=====		

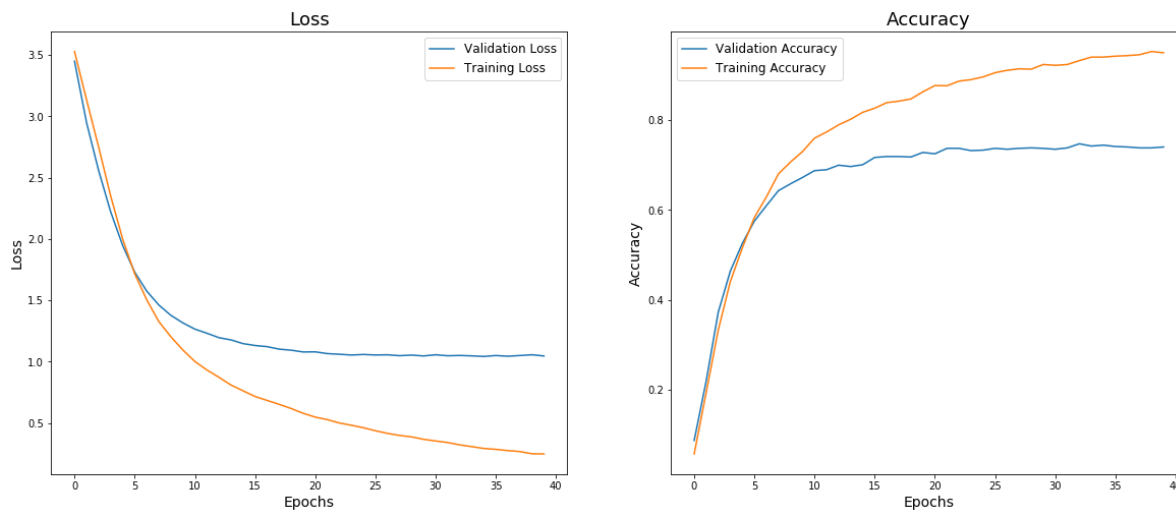
Total params: 192,403

Trainable params: 191,259

Non-trainable params: 1,144

In [29]:

```
1 # Visualise the training and validation loss and accuracy across epochs
2 visualise_results(us_results_3)
```



In [30]:

```
1 # Evaluate the training results
2 us_results_3_train = us_model_2.evaluate(us_X_train, us_train_y)
3 us_results_3_train
```

```
124/124 [=====] - 1s 5ms/step - loss: 0.1908
- acc: 0.9645
```

Out[30]:

```
[0.19082011282444, 0.9644850492477417]
```

In [31]:

```
1 # Evaluate the test results
2 results_4_test = us_model_2.evaluate(us_X_test, us_test_y)
3 results_4_test
```

```
39/39 [=====] - 0s 4ms/step - loss: 1.0853 -
acc: 0.7186
```

Out[31]:

```
[1.085261344909668, 0.7185726165771484]
```

### 4.3.1 Running the model using Mel Spectrograms instead of MFCCs

A thought exercise more than anything else, just to see if there are any performance gains that can be made by using Mel Spectrograms instead of MFCCs.



In [43]:

```
1  # Create the train, test and val datasets for the Ultrasuite top 35 subset using  
2  us_data_path = 'ultrasuite_top35_data_melspec.json'  
3  (  
4      X_train_mel,  
5      y_train_mel,  
6      X_val_mel,  
7      y_val_mel,  
8      X_test_mel,  
9      y_test_mel,  
10 ) = preprocess.create_train_test(us_data_path, 'mel_specs')
```

Datasets loaded...

In [44]:

```
1 # Create model using Mel Spectrograms instead of MFCCs
2 us_mel_input_shape = (X_train_mel.shape[1], X_train_mel.shape[2], 1)
3 us_mel_output_units = 35
4 us_mel_model = build_model_2(us_mel_input_shape,
5                               us_mel_output_units,
6                               learning_rate=0.0001)
7
8 # One-hot encode Speech Commands labels
9 mel_train_y = reformat_y(y_train_mel)
10 mel_test_y = reformat_y(y_test_mel)
11 mel_val_y = reformat_y(y_val_mel)
12
13 # Fit model
14 us_mel_epochs = 40
15 us_mel_batch_size = 64
16 us_mel_patience = 3
17 us_mel_results = fit_model(us_mel_model,
18                             us_mel_epochs,
19                             us_mel_batch_size,
20                             us_mel_patience,
21                             X_train_mel,
22                             mel_train_y,
23                             X_val_mel,
24                             mel_val_y)
```

```
- acc: 0.8311 - val_loss: 3.5039 - val_acc: 0.2779
```

Epoch 32/40

```
62/62 [=====] - 3s 45ms/step - loss: 0.6858
```

```
- acc: 0.8242 - val_loss: 3.4821 - val_acc: 0.2921
```

Epoch 33/40

```
62/62 [=====] - 2s 38ms/step - loss: 0.6960
```

```
- acc: 0.8245 - val_loss: 3.5068 - val_acc: 0.2941
```

Epoch 34/40

```
62/62 [=====] - 2s 39ms/step - loss: 0.6607
```

```
- acc: 0.8351 - val_loss: 3.5081 - val_acc: 0.2779
```

Epoch 35/40

```
62/62 [=====] - 2s 38ms/step - loss: 0.6575
```

```
- acc: 0.8341 - val_loss: 3.5218 - val_acc: 0.2911
```

Epoch 36/40

```
62/62 [=====] - 2s 35ms/step - loss: 0.6337
```

```
- acc: 0.8412 - val_loss: 3.5604 - val_acc: 0.3053
```

Epoch 37/40

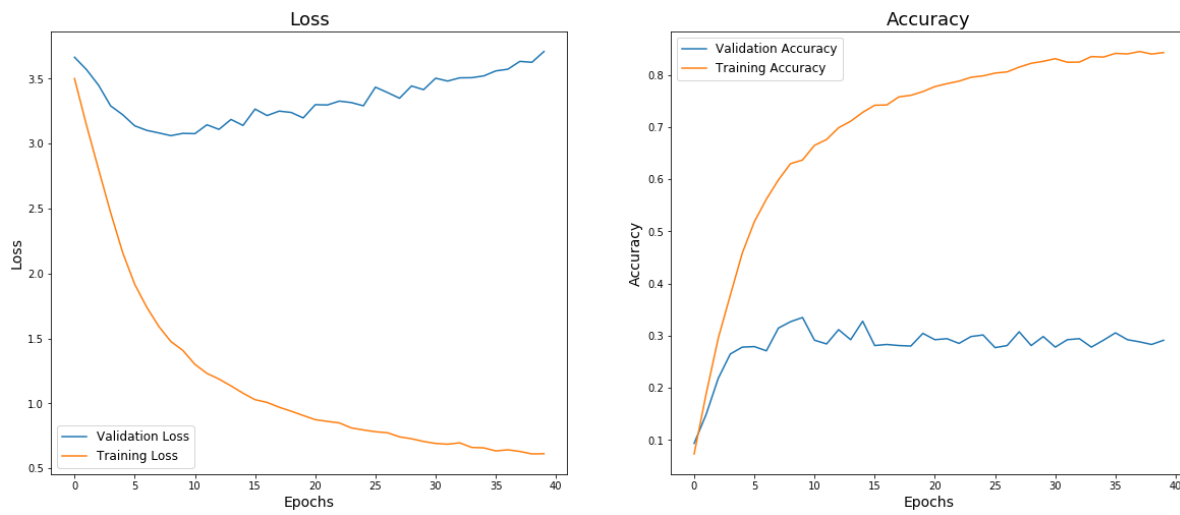
```
62/62 [=====] - 3s 42ms/step - loss: 0.6426
```

```
- acc: 0.8402 - val_loss: 3.5733 - val_acc: 0.2921
```

Epoch 38/40

In [45]:

```
1 # Visualise the training and validation loss and accuracy across epochs
2 visualise_results(us_mel_results)
```



## 4.4 Model 4: Convolutional Neural Network

In [32]:

```
1 # Build a CNN model
2
3 def build_cnn_model(input_shape,
4                     output_units,
5                     loss='categorical_crossentropy',
6                     learning_rate=0.0001):
7
8     cnn_model = tf.keras.models.Sequential()
9
10    # 1st convolutional layer
11    cnn_model.add(tf.keras.layers.Conv2D(32, (3, 3),
12                                         activation='relu',
13                                         input_shape=input_shape,
14                                         kernel_regularizer=tf.keras.regularizers.L2(0.01)))
15    cnn_model.add(tf.keras.layers.BatchNormalization())
16    cnn_model.add(tf.keras.layers.MaxPooling2D((2, 2),
17                                                strides=(2, 2),
18                                                padding='same'))
19
20    # 2nd convolutional layer
21    cnn_model.add(tf.keras.layers.Conv2D(32, (4, 4),
22                                         activation='relu',
23                                         kernel_regularizer=tf.keras.regularizers.L2(0.01)))
24    cnn_model.add(tf.keras.layers.BatchNormalization())
25    cnn_model.add(tf.keras.layers.MaxPooling2D((3, 3),
26                                                strides=(2, 2),
27                                                padding='same'))
28
29    # 3rd convolutional layer
30    cnn_model.add(tf.keras.layers.Conv2D(64, (2, 2),
31                                         activation='relu',
32                                         kernel_regularizer=tf.keras.regularizers.L2(0.01)))
33    cnn_model.add(tf.keras.layers.BatchNormalization())
34    cnn_model.add(tf.keras.layers.MaxPooling2D((2, 2),
35                                                strides=(2, 2),
36                                                padding='same'))
37
38    # Flatten output and feed into dense layer
39    cnn_model.add(tf.keras.layers.Flatten())
40    cnn_model.add(tf.keras.layers.Dense(32, activation='relu'))
41    cnn_model.add(tf.keras.layers.Dropout(0.3))
42
43    # Softmax output layer
44    cnn_model.add(tf.keras.layers.Dense(output_units, activation='softmax'))
45
46    optimiser = tf.optimizers.Adam(learning_rate=learning_rate)
47
48    # Compile model
49    cnn_model.compile(optimizer=optimiser,
50                    loss=loss,
51                    metrics=['acc'])
52
53    # Print summary for model
54    cnn_model.summary()
55
56    return cnn_model
```

## 4.4.1 CNN model for the Speech Commands dataset

In [33]:

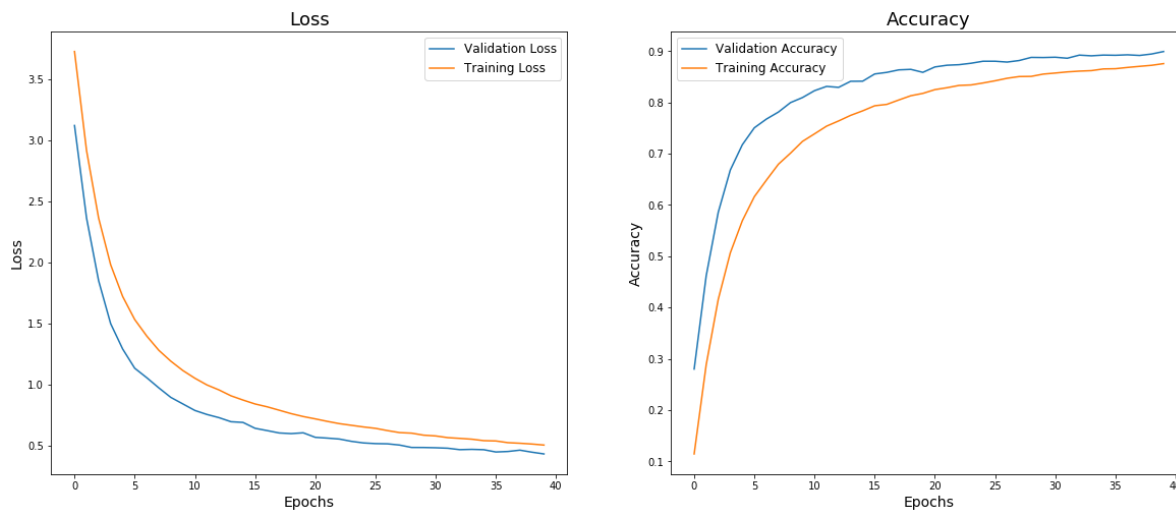
```
1 # Create CNN model using Speech Commands MFCCs
2 sc_cnn_input_shape = (sc_X_train.shape[1], sc_X_train.shape[2], 1)
3 sc_output_units = 35
4 sc_cnn_model = build_cnn_model(sc_cnn_input_shape,
5                               sc_output_units,
6                               learning_rate=0.0001)
7
8
9 # Fit model to Speech Commands data
10 sc_epochs = 40
11 sc_batch_size = 64
12 sc_patience = 3
13 sc_cnn_results = fit_model(sc_cnn_model,
14                            sc_epochs,
15                            sc_batch_size,
16                            sc_patience,
17                            sc_X_train,
18                            sc_train_y,
19                            sc_X_val,
20                            sc_val_y)
```

Model: "sequential\_4"

Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 42, 11, 32)	320
batch_normalization_4 (Batch Normalization)	(None, 42, 11, 32)	128
max_pooling2d (MaxPooling2D)	(None, 21, 6, 32)	0
conv2d_1 (Conv2D)	(None, 18, 3, 32)	16416
batch_normalization_5 (Batch Normalization)	(None, 18, 3, 32)	128
max_pooling2d_1 (MaxPooling2D)	(None, 9, 2, 32)	0
conv2d_2 (Conv2D)	(None, 8, 1, 64)	8256
batch_normalization_6 (Batch Normalization)	(None, 8, 1, 64)	256

In [34]:

```
1 # Visualise the training and validation loss and accuracy across epochs
2 visualise_results(sc_cnn_results)
```



In [35]:

```
1 # Evaluate the training results
2 sc_cnn_results_train = sc_cnn_model.evaluate(sc_X_train, sc_train_y)
3 sc_cnn_results_train
```

1908/1908 [=====] - 17s 9ms/step - loss: 0.31  
48 - acc: 0.9334

Out[35]:

```
[0.3147549331188202, 0.9333846569061279]
```

In [36]:

```
1 # Evaluate the test results
2 sc_cnn_results_test = sc_cnn_model.evaluate(sc_X_test, sc_test_y)
3 sc_cnn_results_test
```

597/597 [=====] - 5s 8ms/step - loss: 0.4375  
- acc: 0.8988

Out[36]:

```
[0.43752631545066833, 0.8987892270088196]
```

## 4.4.2 CNN model for the Ultrasuite dataset

In [37]:

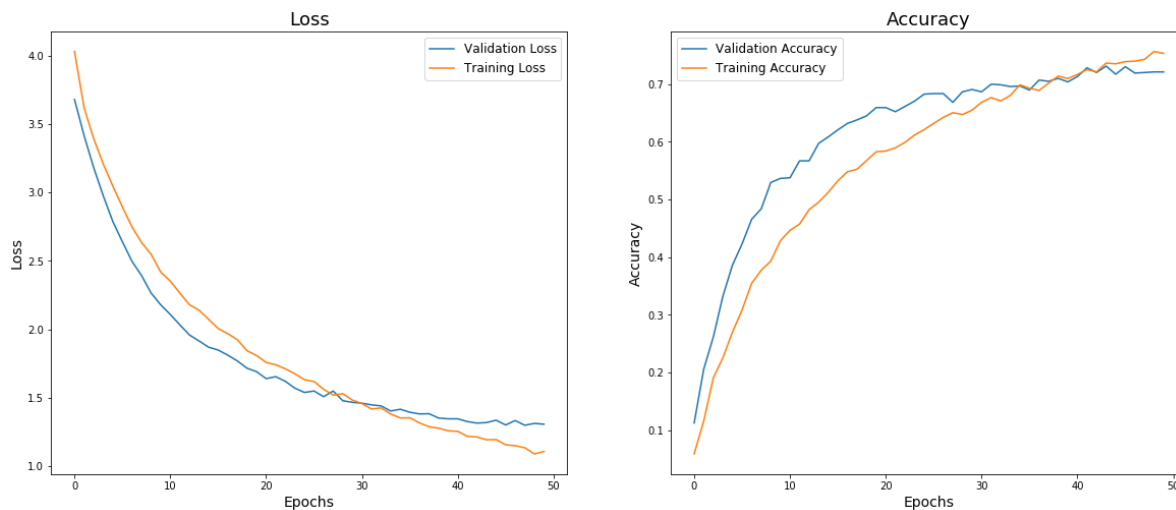
```
1 # Create CNN model using Ultrasuite MFCCs
2 us_cnn_input_shape = (us_X_train.shape[1], us_X_train.shape[2], 1)
3 us_output_units = 35
4 us_cnn_model = build_cnn_model(us_cnn_input_shape,
5                               us_output_units,
6                               learning_rate=0.0001)
7
8
9 # Fit model to Ultrasuite data
10 us_epochs = 50
11 us_batch_size = 16
12 us_patience = 3
13 us_cnn_results = fit_model(us_cnn_model,
14                             us_epochs,
15                             us_batch_size,
16                             us_patience,
17                             us_X_train,
18                             us_train_y,
19                             us_X_val,
20                             us_val_y)
```

Model: "sequential\_5"

Layer (type)	Output Shape	Param #
=====		
conv2d_3 (Conv2D)	(None, 42, 11, 32)	320
<hr/>		
batch_normalization_7 (Batch Normalization)	(None, 42, 11, 32)	128
<hr/>		
max_pooling2d_3 (MaxPooling2D)	(None, 21, 6, 32)	0
<hr/>		
conv2d_4 (Conv2D)	(None, 18, 3, 32)	16416
<hr/>		
batch_normalization_8 (Batch Normalization)	(None, 18, 3, 32)	128
<hr/>		
max_pooling2d_4 (MaxPooling2D)	(None, 9, 2, 32)	0
<hr/>		
conv2d_5 (Conv2D)	(None, 8, 1, 64)	8256
<hr/>		
batch_normalization_9 (Batch Normalization)	(None, 8, 1, 64)	256

In [38]:

```
1 # Visualise the training and validation loss and accuracy across epochs
2 visualise_results(us_cnn_results)
```



In [39]:

```
1 # Evaluate the training results
2 us_cnn_results_train = us_cnn_model.evaluate(us_X_train, us_train_y)
3 us_cnn_results_train
```

```
124/124 [=====] - 1s 10ms/step - loss: 0.7081
- acc: 0.8950: 0s - loss: 0.6985 -
```

Out[39]:

```
[0.70814049243927, 0.8949771523475647]
```

In [40]:

```
1 # Evaluate the test results
2 us_cnn_results_test = us_cnn_model.evaluate(us_X_test, us_test_y)
3 us_cnn_results_test
```

```
39/39 [=====] - 0s 9ms/step - loss: 1.3564 -
acc: 0.7072
```

Out[40]:

```
[1.3563611507415771, 0.7072181701660156]
```

In [41]:

```
1 save_model(us_cnn_model, 'final_model.h5')
```

## 4.5 Final Model Performance Evaluation



In [96]:

```
1 # Import necessary libraries for performance evaluation.
2 from sklearn.metrics import accuracy_score, confusion_matrix
3
4 # Create predictions
5 preds = us_cnn_model.predict(us_X_test)
6
7 # Calculate accuracy and confusion matrix
8 acc = accuracy_score(us_test_y, np.round(preds))*100
9 cm = confusion_matrix(us_test_y.argmax(axis=1),
10                      np.round(preds.argmax(axis=1)))
11
12 print('CONFUSION MATRIX -----')
13 print(cm)
14
15 print('\nTRAIN METRICS -----')
16 print('Loss: {}'.format(us_cnn_results_train[0]))
17 print('Accuracy: {}'.format(np.round(us_cnn_results_train[1]*100), 2))
18
19 print('\nTEST METRICS -----')
20 print('Loss: {}'.format(us_cnn_results_test[0]))
21 print('Accuracy: {}'.format(np.round(us_cnn_results_test[1]*100), 2))
```

CONFUSION MATRIX -----

```
[[14  1  0 ...  1  1  0]
 [ 0 35  0 ...  0  3  0]
 [ 0  0 23 ...  0  0  0]
 ...
 [ 0  0  1 ... 16  0  0]
 [ 1  0  0 ...  0 45  0]
 [ 0  0  0 ...  0  0 14]]
```

TRAIN METRICS -----

Loss: 0.70814049243927  
Accuracy: 89.0%

TEST METRICS -----

Loss: 1.3563611507415771  
Accuracy: 71.0%

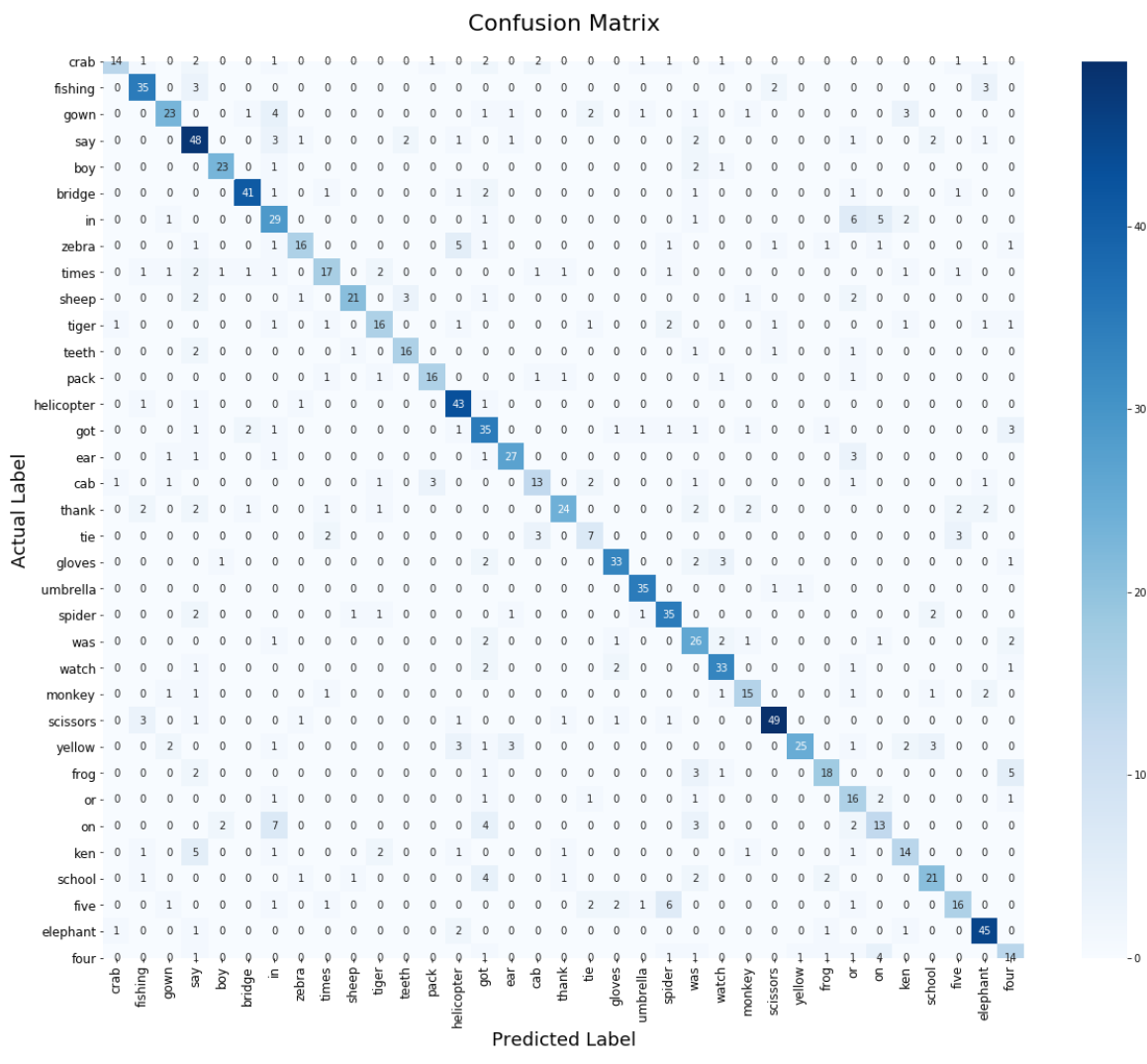
In [92]:

```
1 # Import necessary libraries for classification report
2 from sklearn.metrics import classification_report
3
4 # Print the classification report
5 print(classification_report(us_test_y.argmax(axis=1),
6                             preds.argmax(axis=1),
7                             target_names=us_keywords))
```

	precision	recall	f1-score	support
crab	0.82	0.50	0.62	28
fishing	0.78	0.81	0.80	43
gown	0.74	0.61	0.67	38
say	0.61	0.77	0.68	62
boy	0.85	0.85	0.85	27
bridge	0.89	0.84	0.86	49
in	0.52	0.64	0.57	45
zebra	0.76	0.55	0.64	29
times	0.68	0.55	0.61	31
sheep	0.88	0.68	0.76	31
tiger	0.67	0.59	0.63	27
teeth	0.76	0.73	0.74	22
pack	0.80	0.73	0.76	22
helicopter	0.73	0.91	0.81	47
got	0.56	0.71	0.63	49
ear	0.82	0.79	0.81	34
cab	0.65	0.54	0.59	24
thank	0.83	0.62	0.71	39
tie	0.47	0.47	0.47	15
gloves	0.82	0.79	0.80	42
umbrella	0.88	0.95	0.91	37
spider	0.71	0.81	0.76	43
was	0.52	0.72	0.60	36
watch	0.77	0.82	0.80	40
monkey	0.68	0.65	0.67	23
scissors	0.89	0.84	0.87	58
yellow	0.93	0.61	0.74	41
frog	0.75	0.60	0.67	30
or	0.40	0.70	0.51	23
on	0.50	0.42	0.46	31
ken	0.58	0.52	0.55	27
school	0.72	0.64	0.68	33
five	0.67	0.52	0.58	31
elephant	0.80	0.88	0.84	51
four	0.48	0.56	0.52	25
accuracy			0.71	1233
macro avg	0.71	0.68	0.69	1233
weighted avg	0.72	0.71	0.71	1233

In [77]:

```
1  # Use Seaborn to make the confusion matrix more visually presentable
2  plt.figure(figsize=(20,16))
3  ax = plt.subplot()
4  sns.heatmap(cm, annot=True, ax=ax, fmt='g', cmap='Blues')
5
6  us_keywords = [
7      'crab',
8      'fishing',
9      'gown',
10     'say',
11     'boy',
12     'bridge',
13     'in',
14     'zebra',
15     'times',
16     'sheep',
17     'tiger',
18     'teeth',
19     'pack',
20     'helicopter',
21     'got',
22     'ear',
23     'cab',
24     'thank',
25     'tie',
26     'gloves',
27     'umbrella',
28     'spider',
29     'was',
30     'watch',
31     'monkey',
32     'scissors',
33     'yellow',
34     'frog',
35     'or',
36     'on',
37     'ken',
38     'school',
39     'five',
40     'elephant',
41     'four']
42
43  ax.set_title('Confusion Matrix', fontsize=22, pad=30)
44  ax.set_xlabel('Predicted Label', fontsize=18)
45  ax.set_ylabel('Actual Label', fontsize=18)
46  ax.xaxis.set_ticklabels(us_keywords, rotation=90, fontsize=12)
47  ax.yaxis.set_ticklabels(us_keywords, rotation=0, fontsize=12)
48  plt.show();
```



## 4.5.1 Making a prediction on an unseen audio sample

In [48]:

```
1 # Load the audio sample and preview
2 target_sample = 'audio/martha-frog.wav'
3 target_label = 'Frog'
4 audio_sample, sr = librosa.load(target_sample)
5 print('Audio sample:', target_label)
6 ipd.Audio(audio_sample, rate=sr)
```

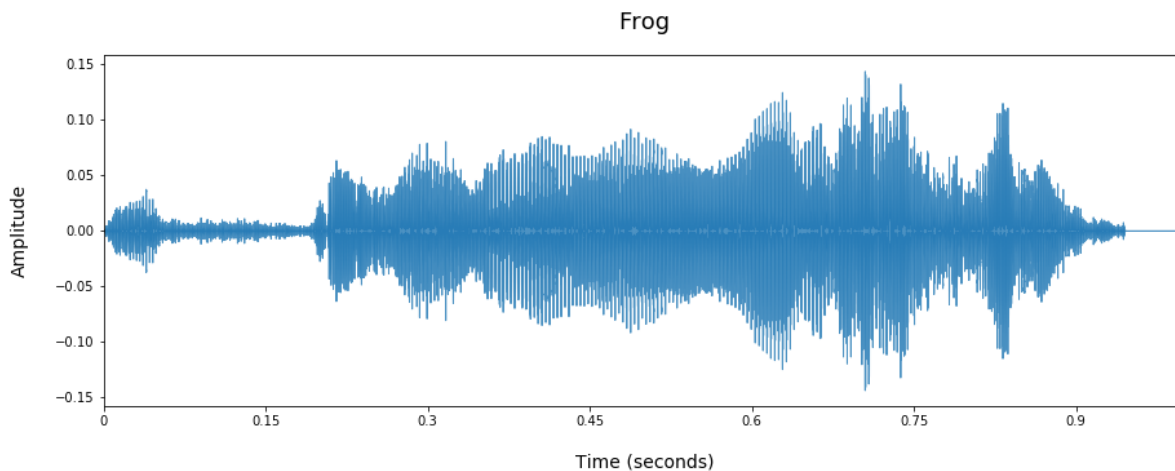
Audio sample: Frog

Out[48]:

▶ 0:00 / 0:01 🔊 ⋮

In [49]:

```
1 # Plot the waveform for the specific audio sample
2 plt.figure(figsize=(15, 5))
3 plt.title(target_label, fontsize=18, pad=20)
4 librosa.display.waveplot(audio_sample, sr, alpha=0.8)
5 plt.xlabel('Time (seconds)', fontsize=14, labelpad=20)
6 plt.ylabel('Amplitude', fontsize=14, labelpad=20)
7 plt.show();
```



In [50]:

```
1 # Run inference on the unseen audio file
2 mfccs = librosa.feature.mfcc(audio_sample,
3                               sr,
4                               n_mfcc=13,
5                               n_fft=2048,
6                               hop_length=512)
7 mfccs = mfccs.T
8 mfccs = mfccs[np.newaxis, ..., np.newaxis]
9
10 prediction = us_cnn_model.predict(mfccs)
11 predicted_index = np.argmax(prediction)
12
13 predicted_keyword = us_keywords[predicted_index]
14 print('Martha says...', predicted_keyword, '!')
```

Martha says... frog !

## 4.6 Overall Conclusion

Initial models did not generalise well and tended to overfit the training data. Subsequent changes to parameters significantly improved the training accuracy but, again were overfitting the training data.

Whilst it does not have a high accuracy score, this final model using a Convolutional Neural Network produced the best results classifying over 70% of the “unseen” audio samples whilst minimising the overfitting to training data. Looking at the classification report, it is also clear that the model performs better for more identifiable words such as `scissors` or `umbrella`.

The very nature of speech sound disorders mean that a model that has been simply trained on audio samples of "typical" speech will generally be more accurate than one that has been trained on audio samples of "atypical" speech as demonstrated above with the model comparison between the Speech Commands and Ultrasuite datasets.

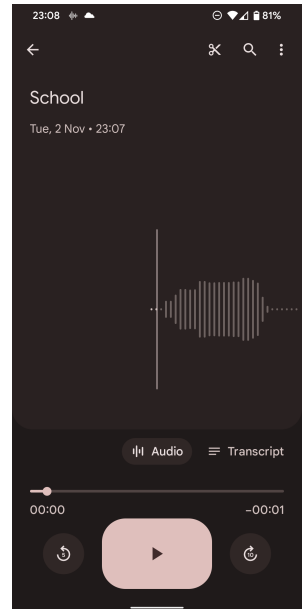
### 4.6.1 Recommendations

It is envisaged that the model would be used in the form of a mobile app, similar to that of Google Recorder (see image), that could be activated by the parent in order to capture the child speaking.

The app would also enable the parent to crop or isolate the audio sample and would provide a side-by-side prediction or transcript of what the child was saying.

This could even be provided as a setting within the Google recorder app itself that allowed the user, in this case the parent, to switch between the standard speech recognition model and our final model specifically trained for children with speech disorders.

Given this eventual usage of this model, it is arguable the app would be more useful if it suggested three potential words in order of likelihood, giving the parent options of what the child might be trying to communicate.



The accuracy of the final model is suitable enough to go ahead with a soft launch of the app to a controlled group of parents, in part for testing but also as a way of collecting additional data to improve the model.

### 4.6.2 Future Work

Future work to improve the model could include:

- Utilise other model architectures and enable them to accept longer audio samples as these are more representative of atypical speech patterns.
- Source additional data in the form of further audio samples potentially even using the app as a means for gathering additional samples and improving the model.
- Use data augmentation when training the model, in particular the [MixSpeech](https://arxiv.org/abs/2102.12664) (<https://arxiv.org/abs/2102.12664>) method that could take a weighted combination of mel-spectrograms and MFCCs in order to improve model performance.

---

#### Sources / Code adapted from:

\* [Hands-On Machine Learning with Scikit-Learn, Keras & Tensorflow - Aurélien Géron](https://www.oreilly.com/library/view/hands-on-machine-learning/9781492032632/)  
(<https://www.oreilly.com/library/view/hands-on-machine-learning/9781492032632/>).

\* [Deep Learning Audio Application from Design to Deployment - Valerio Velardo - The Sound of AI](https://github.com/musikalkemist/Deep-Learning-Audio-Application-From-Design-to-Deployment)  
(<https://github.com/musikalkemist/Deep-Learning-Audio-Application-From-Design-to-Deployment>).

