

By
Clyde Clarke

SPEECH CLASSIFICATION USING CNN, SVM, REGRESSION TREES

Outline

Problem Statement

- Dataset

Sound Classification

Feature Extraction

- MFCC
- STFT
- WAVELET

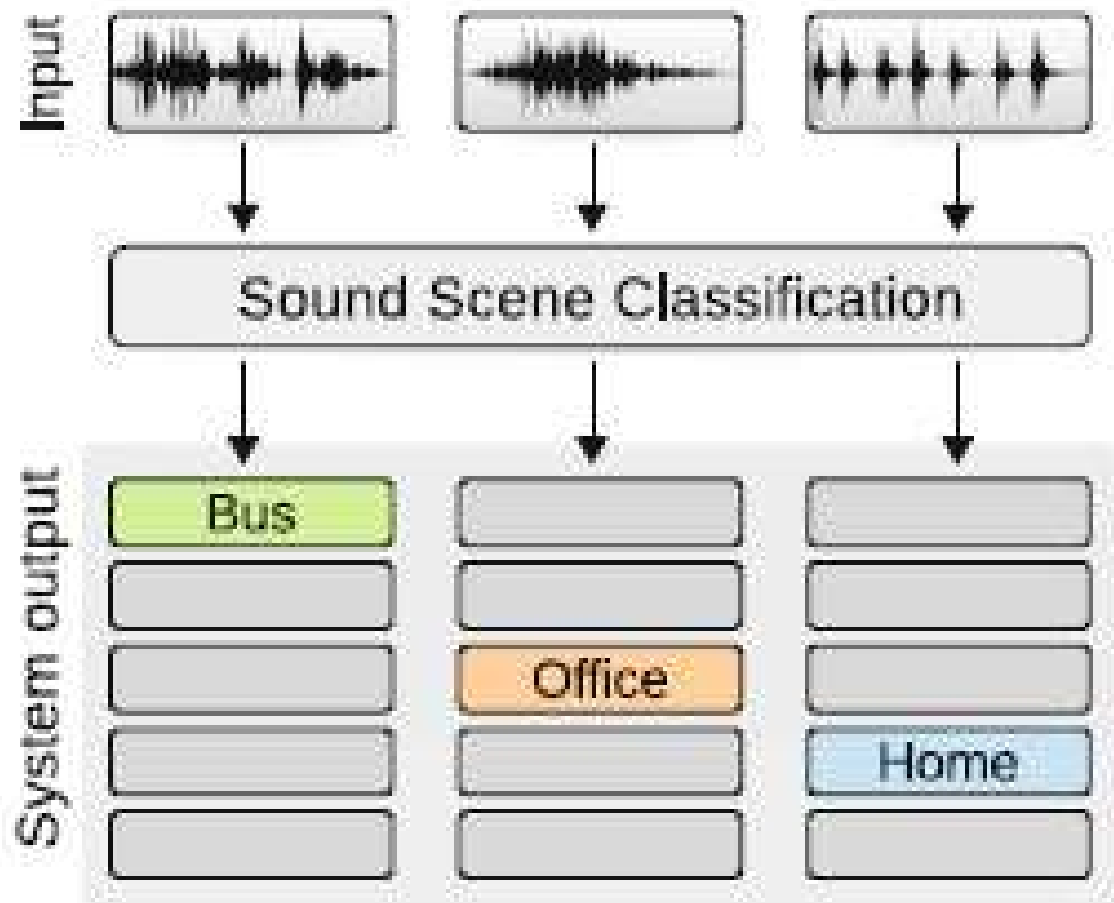
Classification Methods

- SVM
- Regression Trees
- CNN

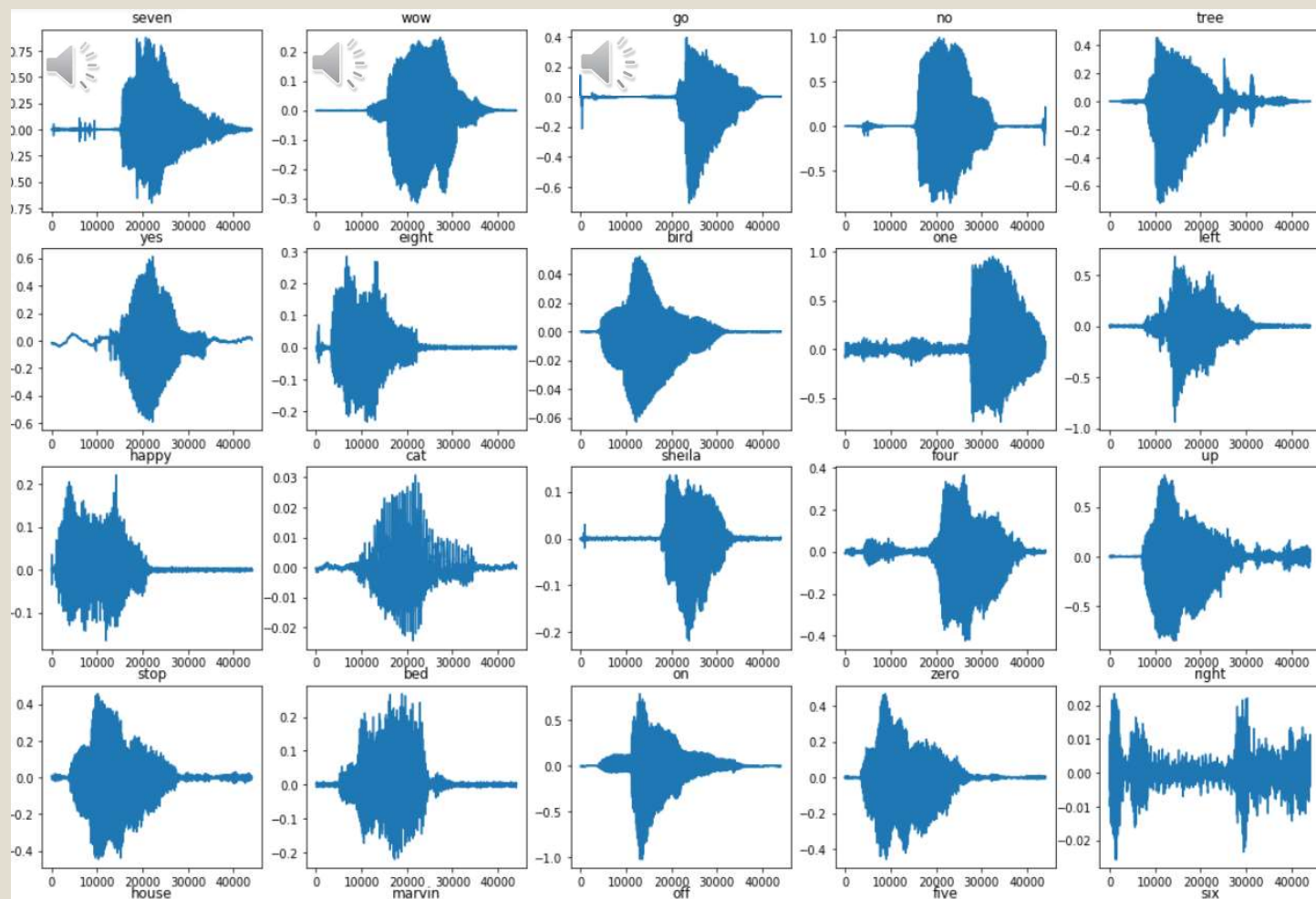
Evaluation

The Speech Dataset

- Speech Commands Dataset
- Contains 15000 labeled sound excerpts (1s) of real field-recording speech from 30 classes : (1).nine, (2).eight, (3).one, (4). zero, (5). bird, (6). shiela, (7). night, (8) siz, (9) stop, up (10.)... etc.
- Sampling Frequency 44 KHz
- Sample Duration ~ 1s
- Sound Samples 15000



Sound
Classification

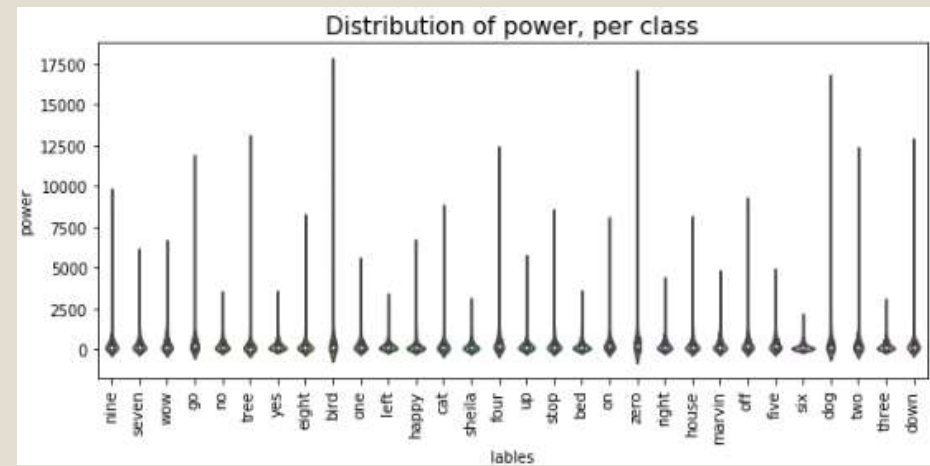
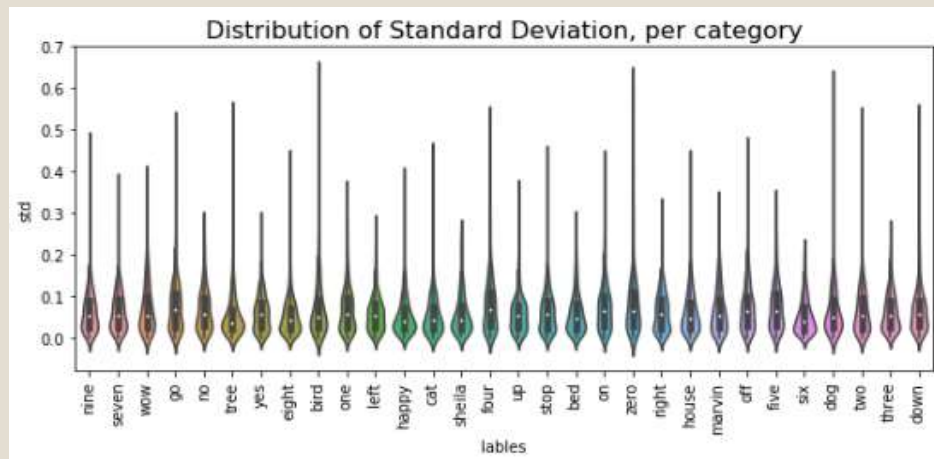
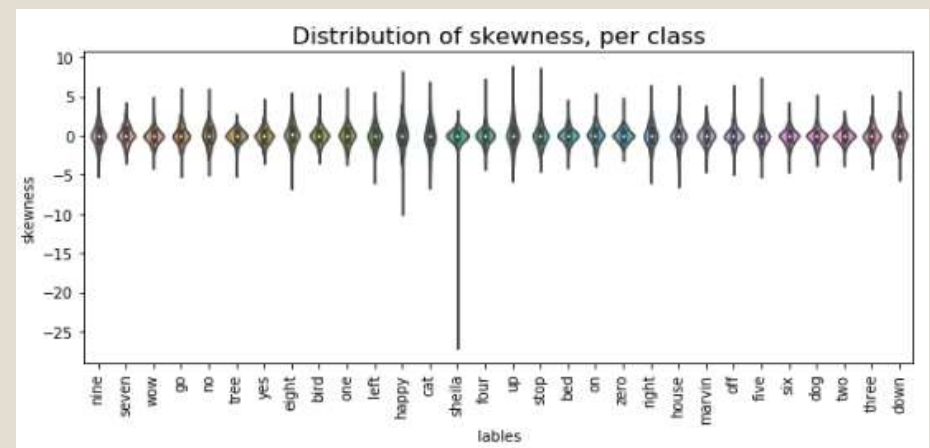
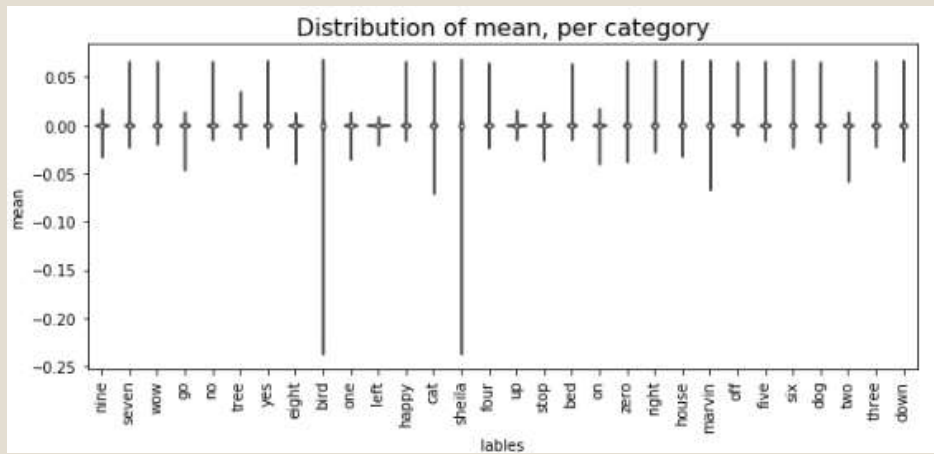


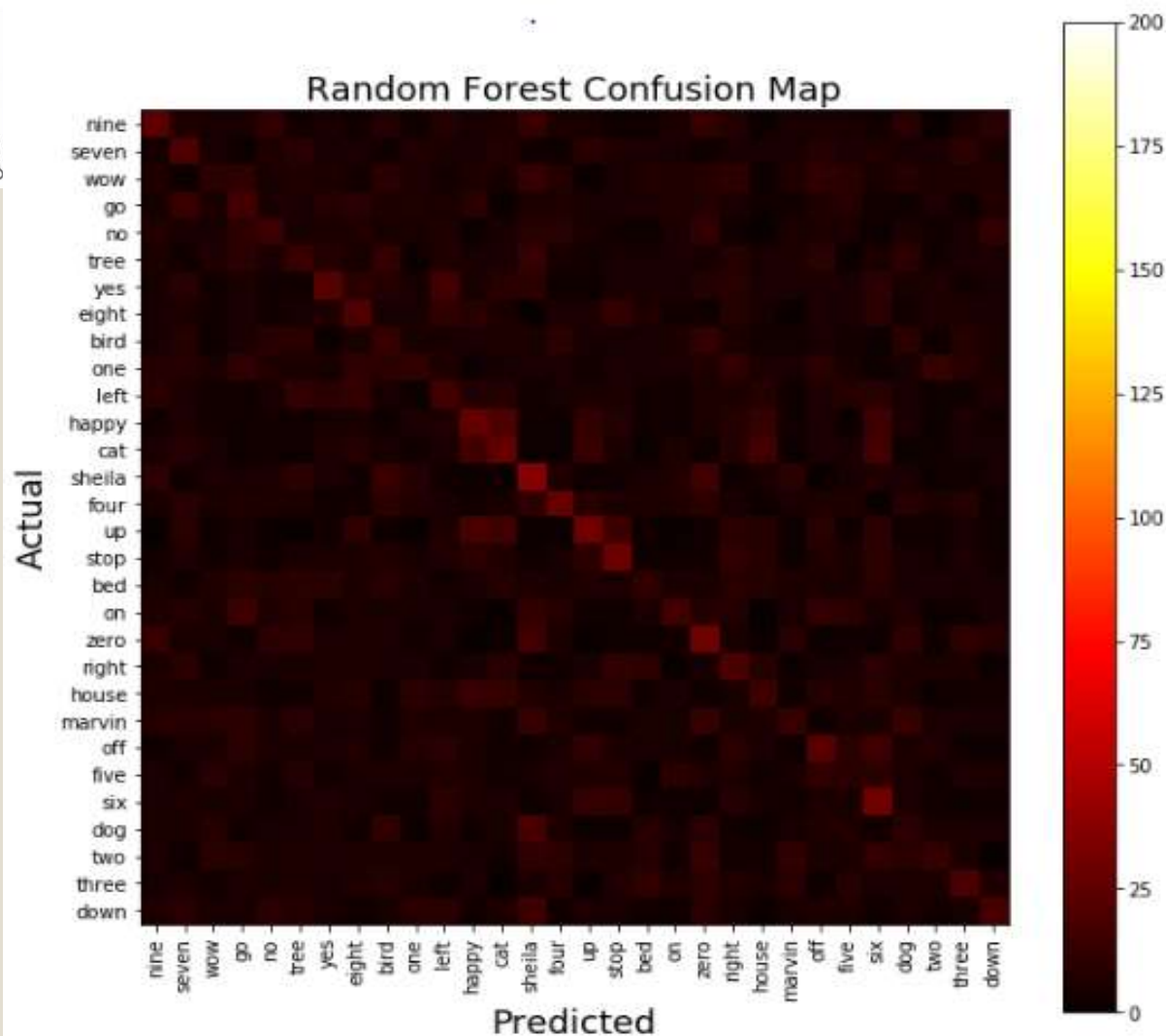
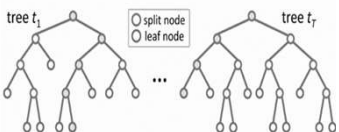
SOUND CLASSES

Sound samples of 1 second in duration

Each audio class is 44100 sample

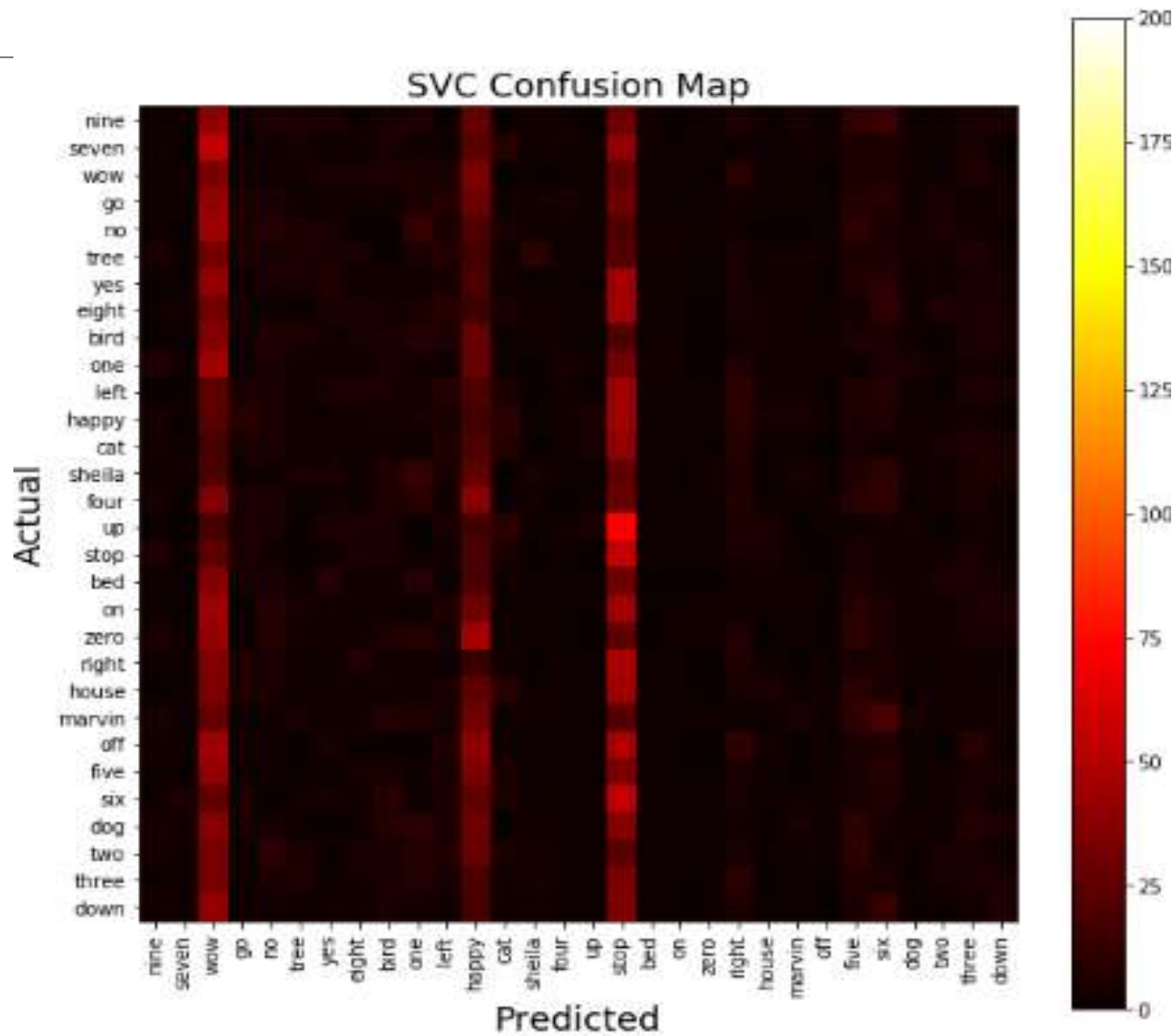
Supervised Learning Features





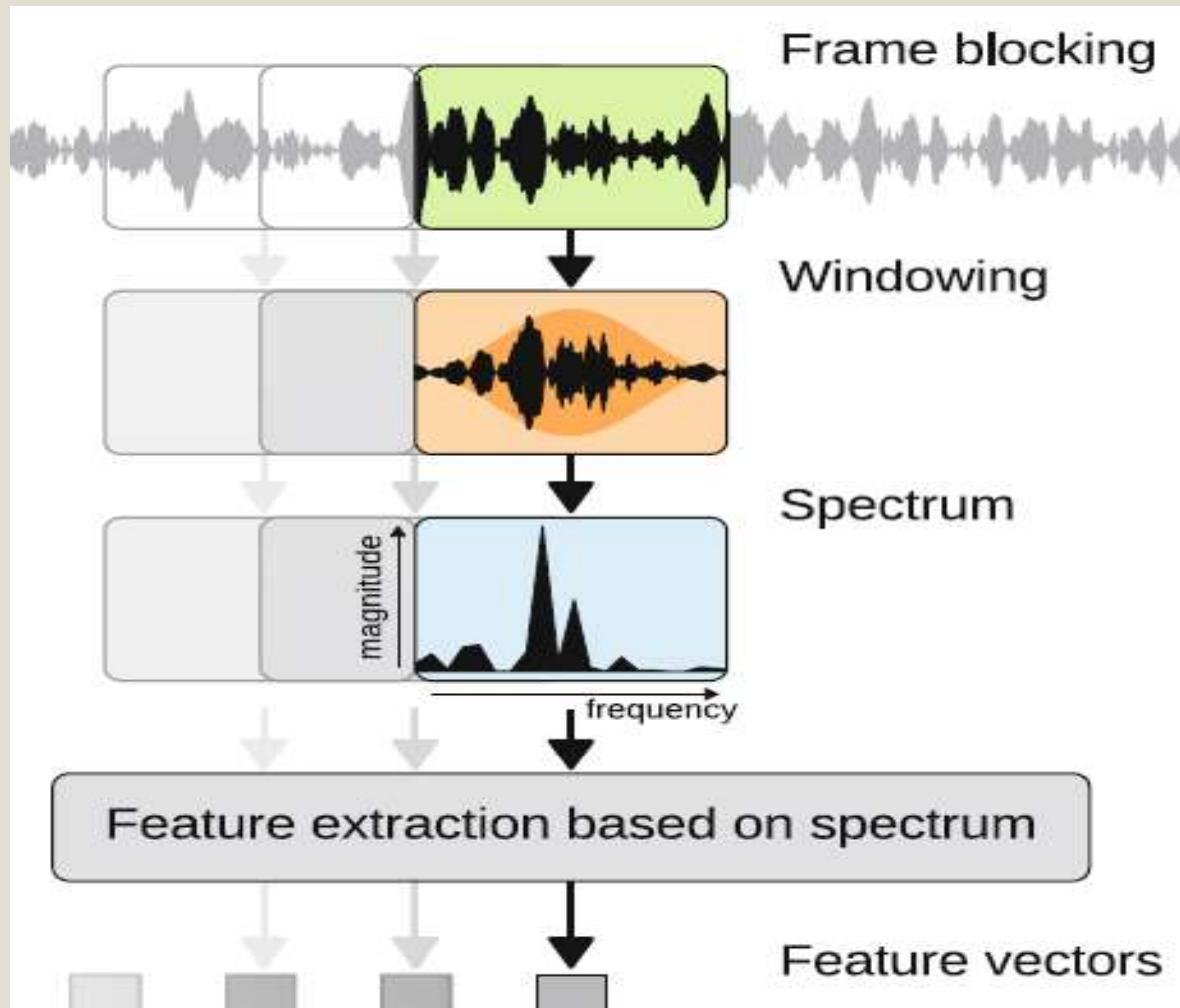
Random Forest

- Random forests (RF) are a combination of tree predictors
- Each tree depends on the values of a random vector sampled in dependently
- The generalization error depends on the strength of the individual trees and the correlation between them
- Accuracy
 - .11



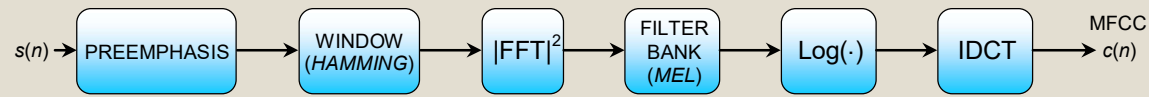
SVM

- Large Margin Classifier
- Maximizes the distance between decision boundary and support vectors
- Kernel
 - Radial Basis Function
- Gamma
 - .08
- Accuracy
 - .11

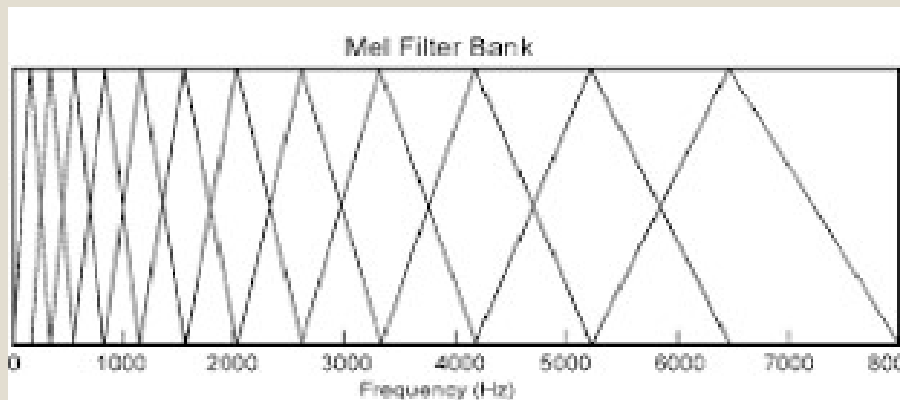


Feature Extraction

- Frame Blocking
 - Segment signal into regions that retain small scale changes
- Windowing
 - Done to prevent gibbs phenomena at frame edges that would corrupt spectral analysis
- Spectrum
 - Any number of techniques can be used to extract features from signal
 - DWT
 - FFT
 - DCT



- Preemphasis: compensates for spectral tilt (speech production/microphone channel)
- Windowing: suppression of transient effects in short-term segments of signal
- $|FFT|^2$: energy spectrum (phase is discarded)
- MEL Filter bank: MEL scale – models logarithmic perception of frequency in humans; triangular filters – dimensionality reduction

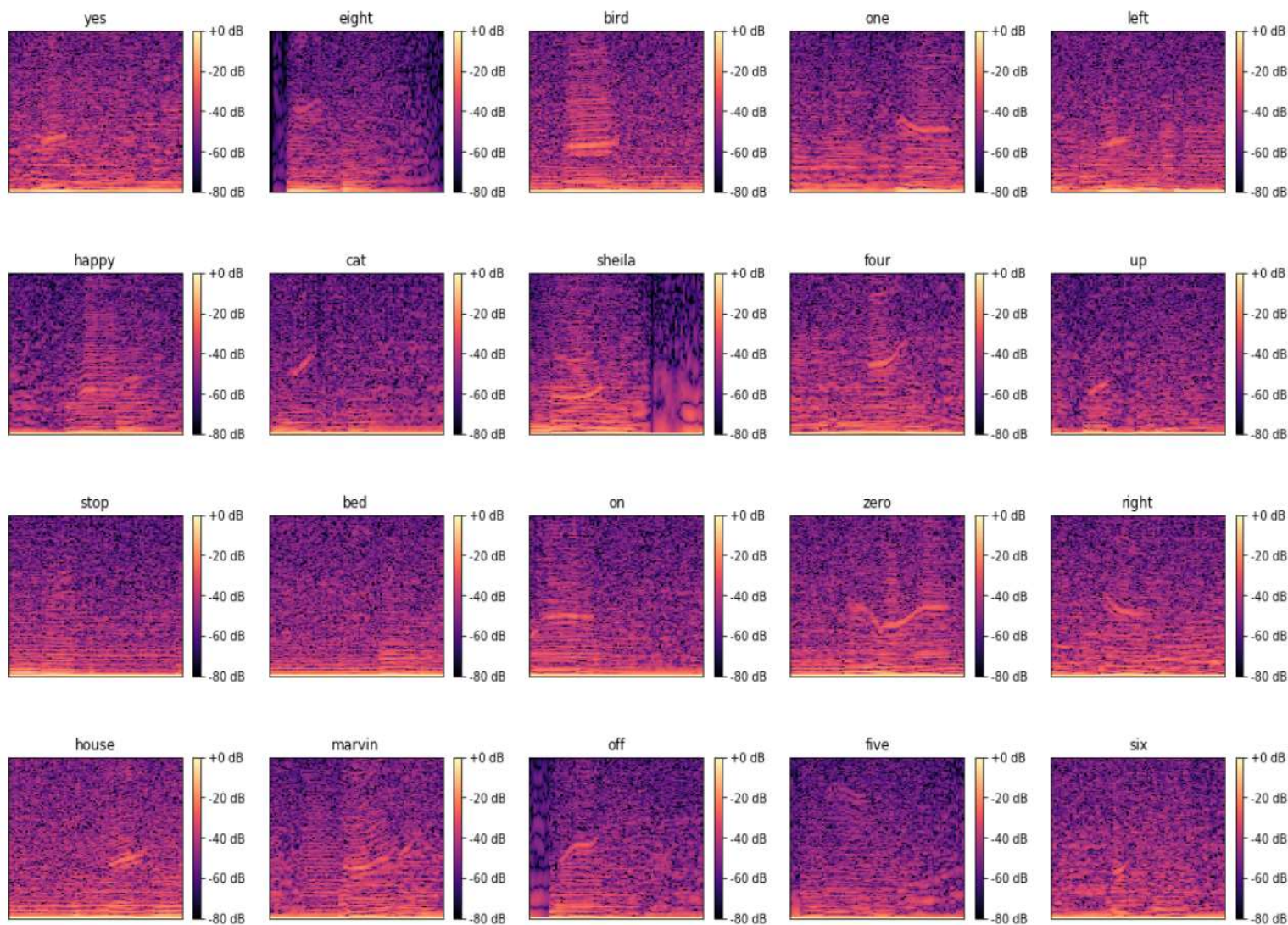


MFCC Features

Pre emphasis filter performs lowpass at Nyquist frequency of signal

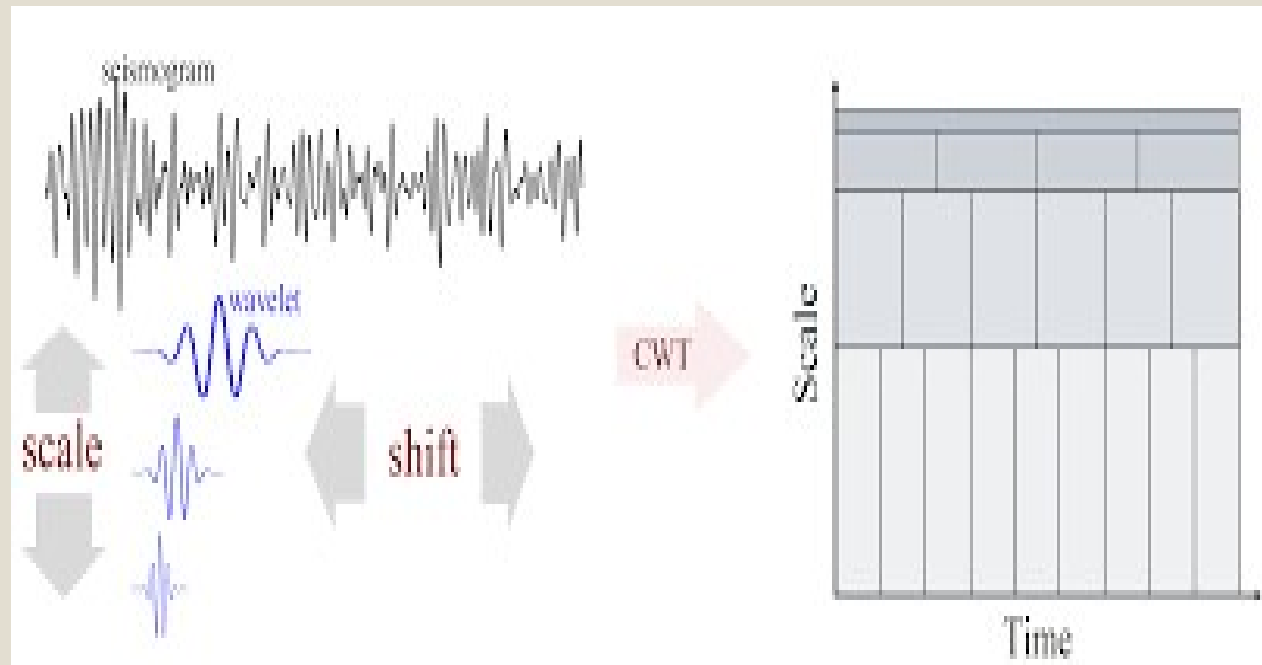
Short Time Fourier transform is performed on the windowed signal

Logrithmic filter bank is used that models the human perception of sound



MFCC FEATURES

Figure shows a
representative sample
from 16 distinct sound
classes



$$(W_a f) = \int f(x) \psi_{a,b}(x) dx$$

with $\psi_{a,b}(x) = \frac{1}{\sqrt{a}} \psi\left(\frac{x-b}{a}\right)$

Wavelet Transform

Continuous Wavelet Transform

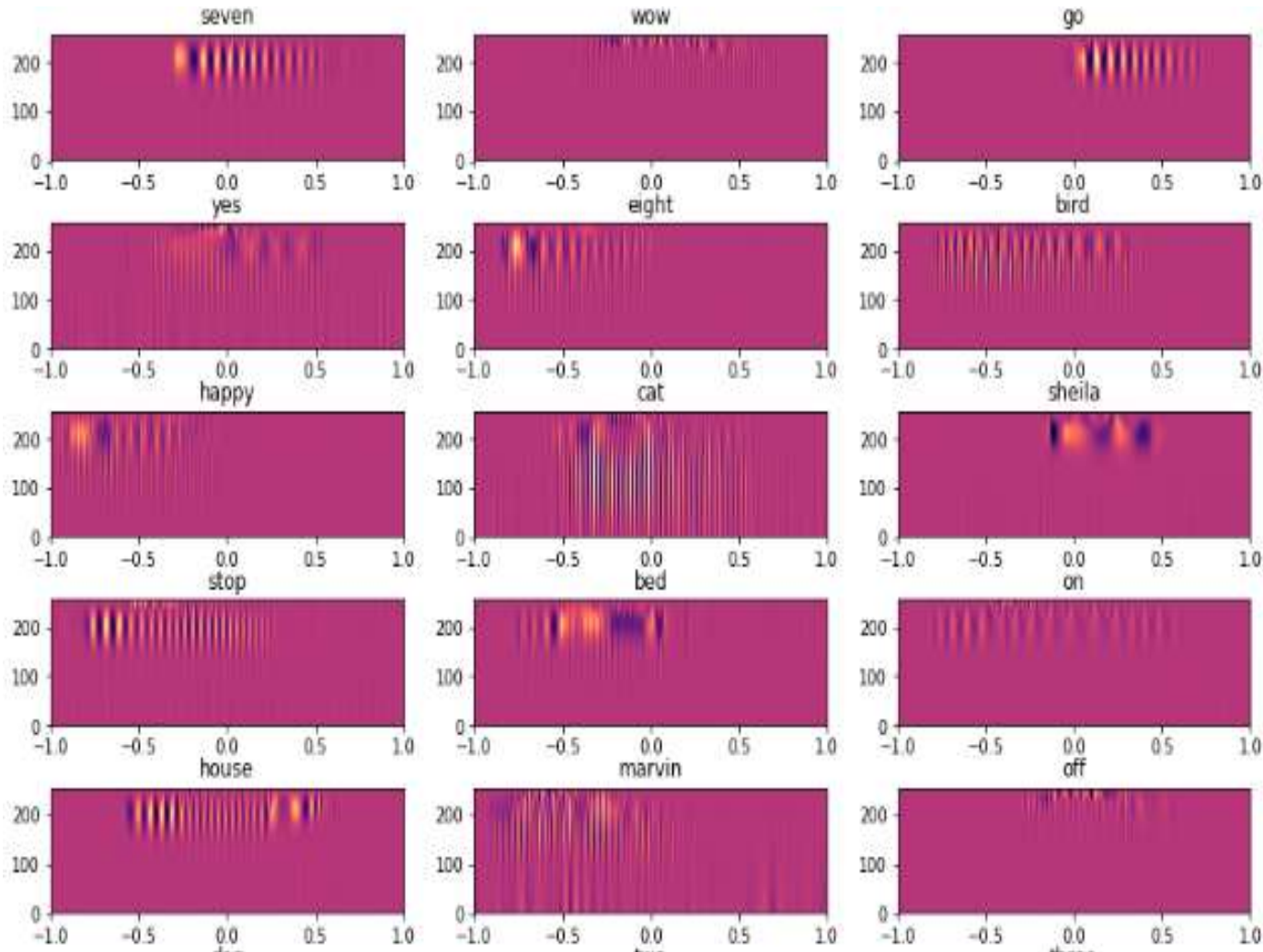
Decomposes signals with (wavelets) orthogonal basis functions

Scalogram is computed by shifting and scaling (ie stretching) wavelet function

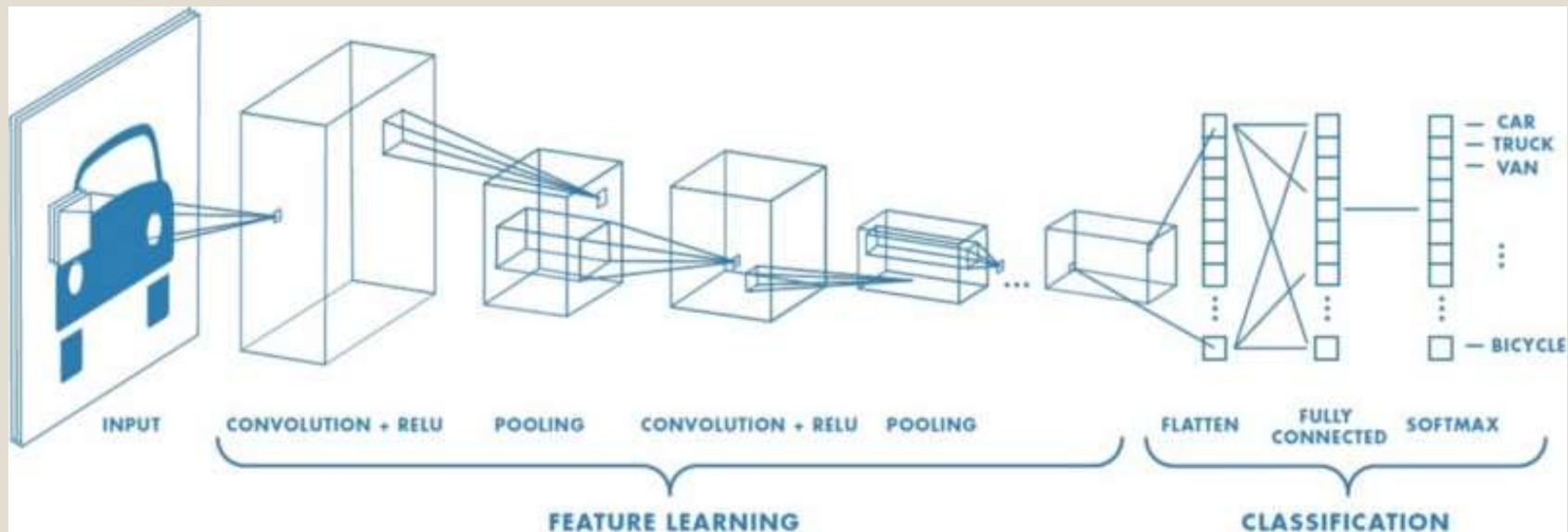
Similar to Fourier Transform

WAVELETS

Continuous Wavelet Transform
(CWT) using Mexican Hat
Wavelet



Convolutional Neural Network Architecture



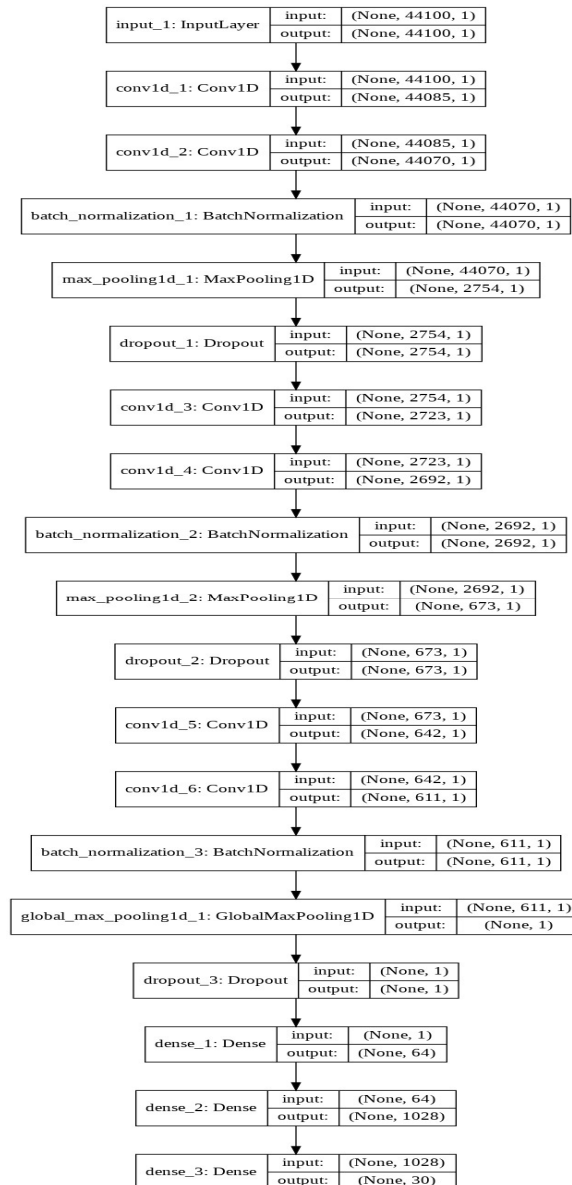
Training and Processing

Due to variable sample length wave files were zero padded

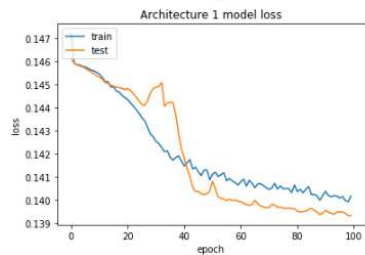
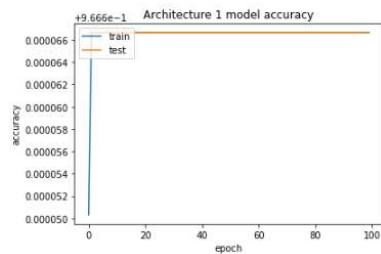
Training was done using an 60/40 split of the data

Network was trained using adam optimization

A batch size of 256 and 100 epoch was used to train five convnet architectures



CNN ARCHITECTURE 1



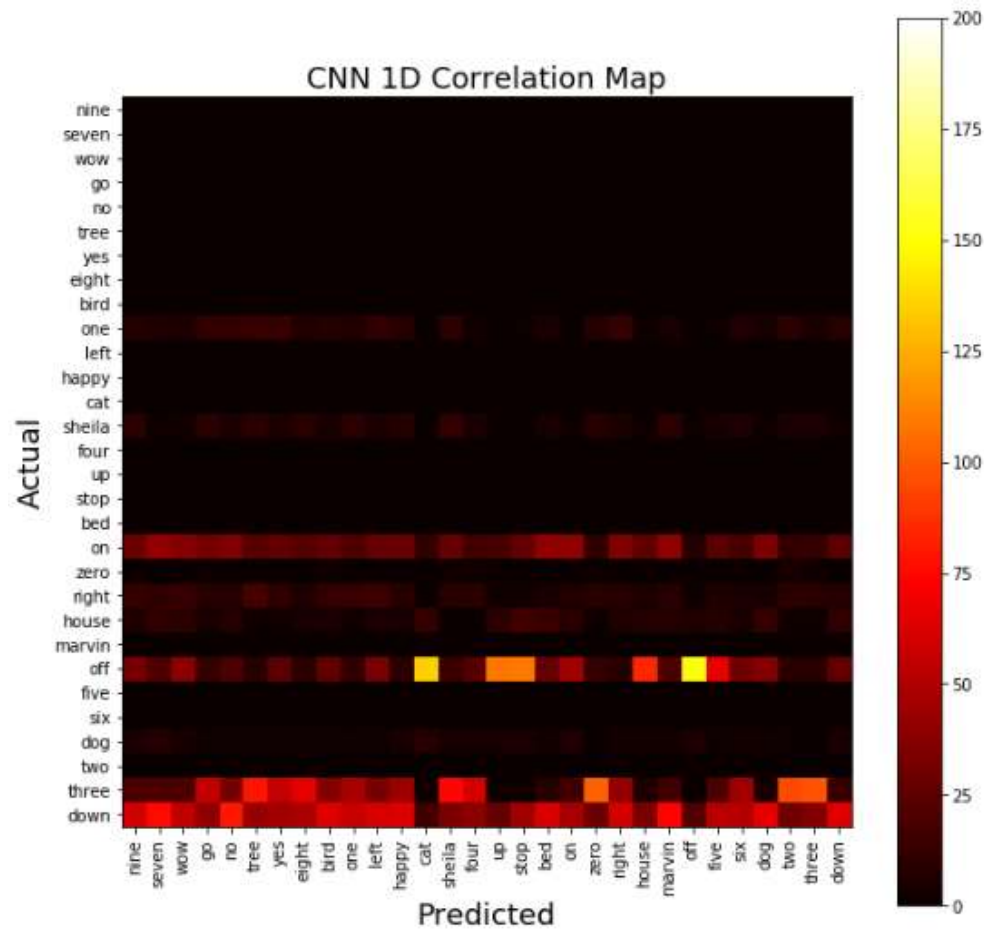
Architecture 1

- Training of the model for both validation and test data shows poor performance
- The classification report also shows poor precision, recall and f-scores across all classes

Classification report

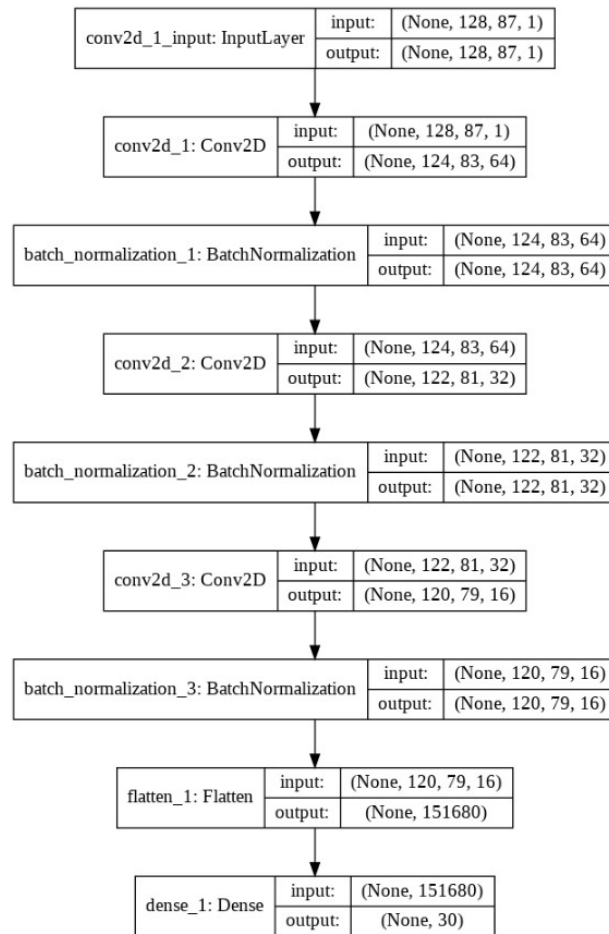
Class	Precision	Recall	F1-score
28 (153)	0.07	0.08	0.08
27 (147)	0.08	0.08	0.08
26 (122)	0.04	0.06	0.05
25 (217)	0.15	0.13	0.14
24 (153)	0.06	0.08	0.07
23 (187)	0.10	0.09	0.09
22 (157)	0.08	0.09	0.09
21 (179)	0.06	0.06	0.06
19 (166)	0.08	0.10	0.09
19 (141)	0.13	0.11	0.12
18 (141)	0.10	0.13	0.11
17 (113)	0.05	0.08	0.06
16 (173)	0.11	0.12	0.11
15 (278)	0.21	0.14	0.17
14 (114)	0.07	0.07	0.07
13 (262)	0.20	0.14	0.17
12 (199)	0.12	0.10	0.11
11 (194)	0.11	0.11	0.11
10 (201)	0.07	0.07	0.07
9 (164)	0.07	0.08	0.08
8 (181)	0.09	0.09	0.09
7 (182)	0.07	0.07	0.07
6 (196)	0.09	0.08	0.08
5 (218)	0.06	0.05	0.05
4 (182)	0.08	0.08	0.08
3 (157)	0.09	0.11	0.10
2 (166)	0.06	0.07	0.07
1 (202)	0.07	0.06	0.06
0 (243)	0.17	0.11	0.13

Color scale: 0.04 to 0.20



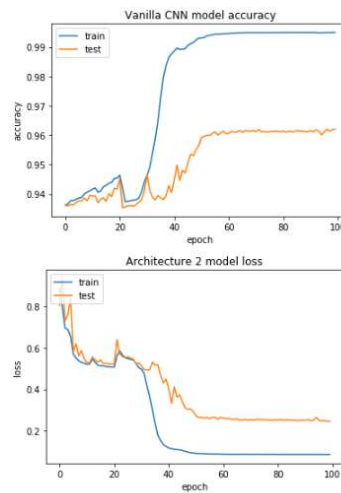
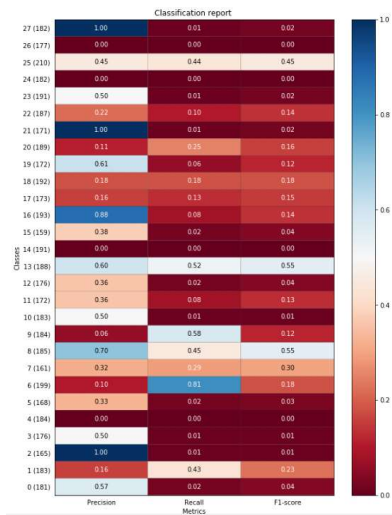
Correlation Map

- The correlation map of the 1D model shows poor performance as it mistakes classes
- The model accuracy is merely 10%



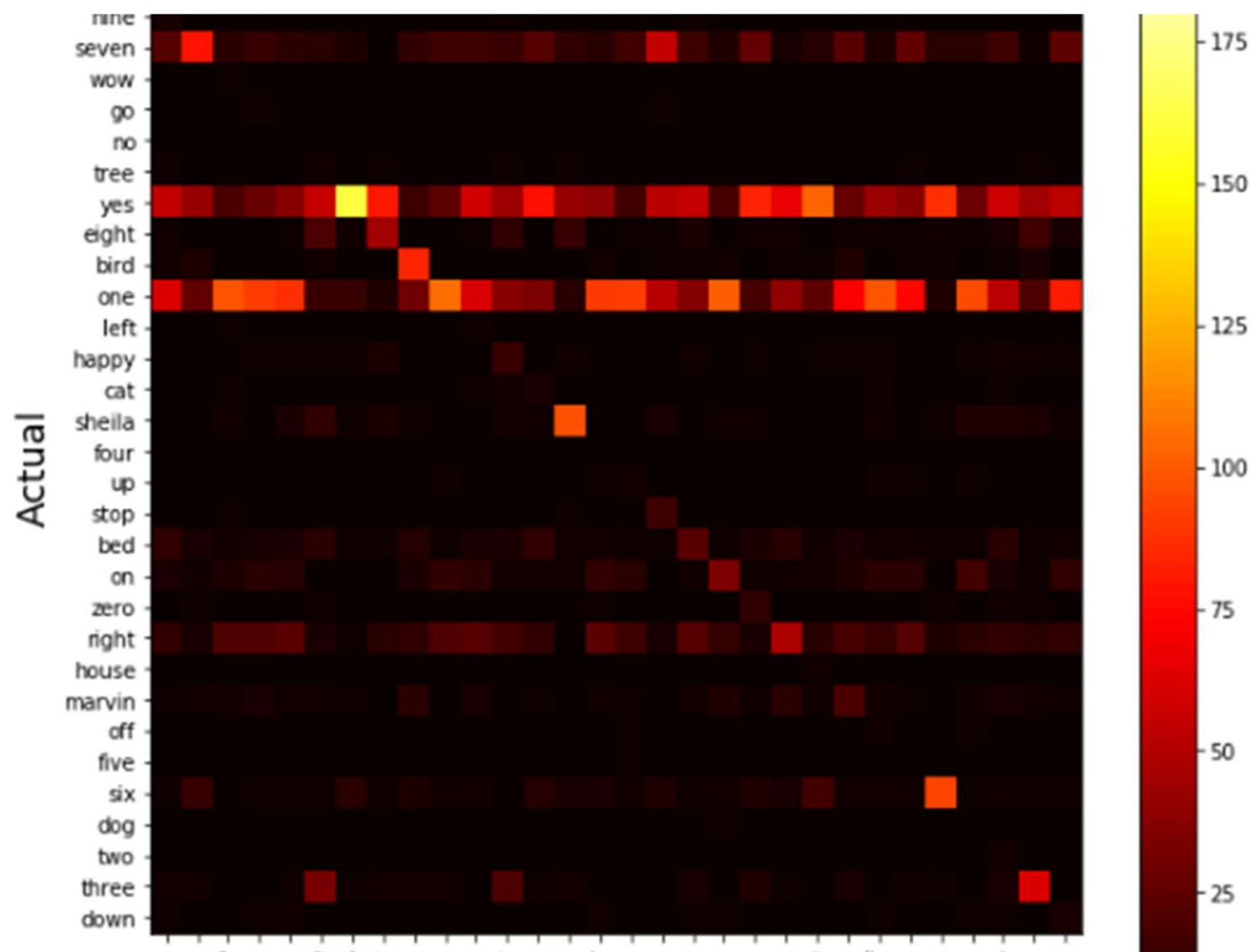
ARCHITECTURE 2

Input features for this
architecture are 2D MFCC
with 200 filters



Performance

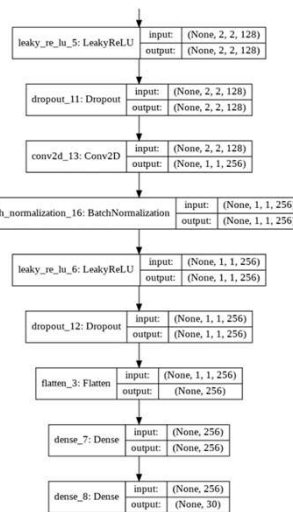
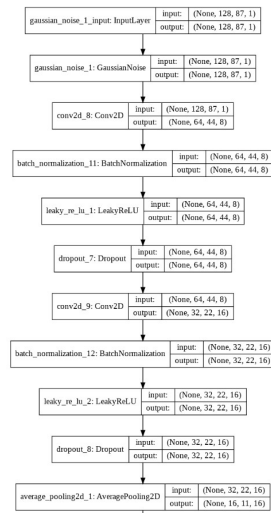
- Trained model has increasing loss and increasing accuracy over 100 epochs
- The performance however is still bad for this model
- Accuracy
 - 12%

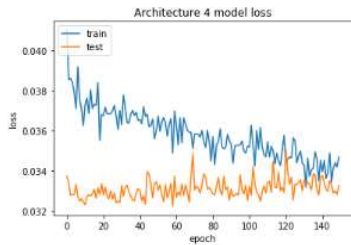
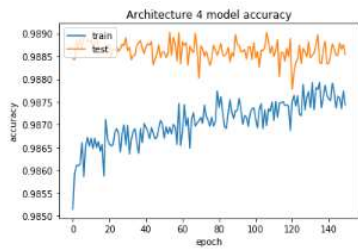


Correlation Map

- Training of this model shows that classification of yes and one had high correlation with all other audio classes in the dataset
- This is probably due to not having enough features ie convolutional layers in the model

Architecture 3

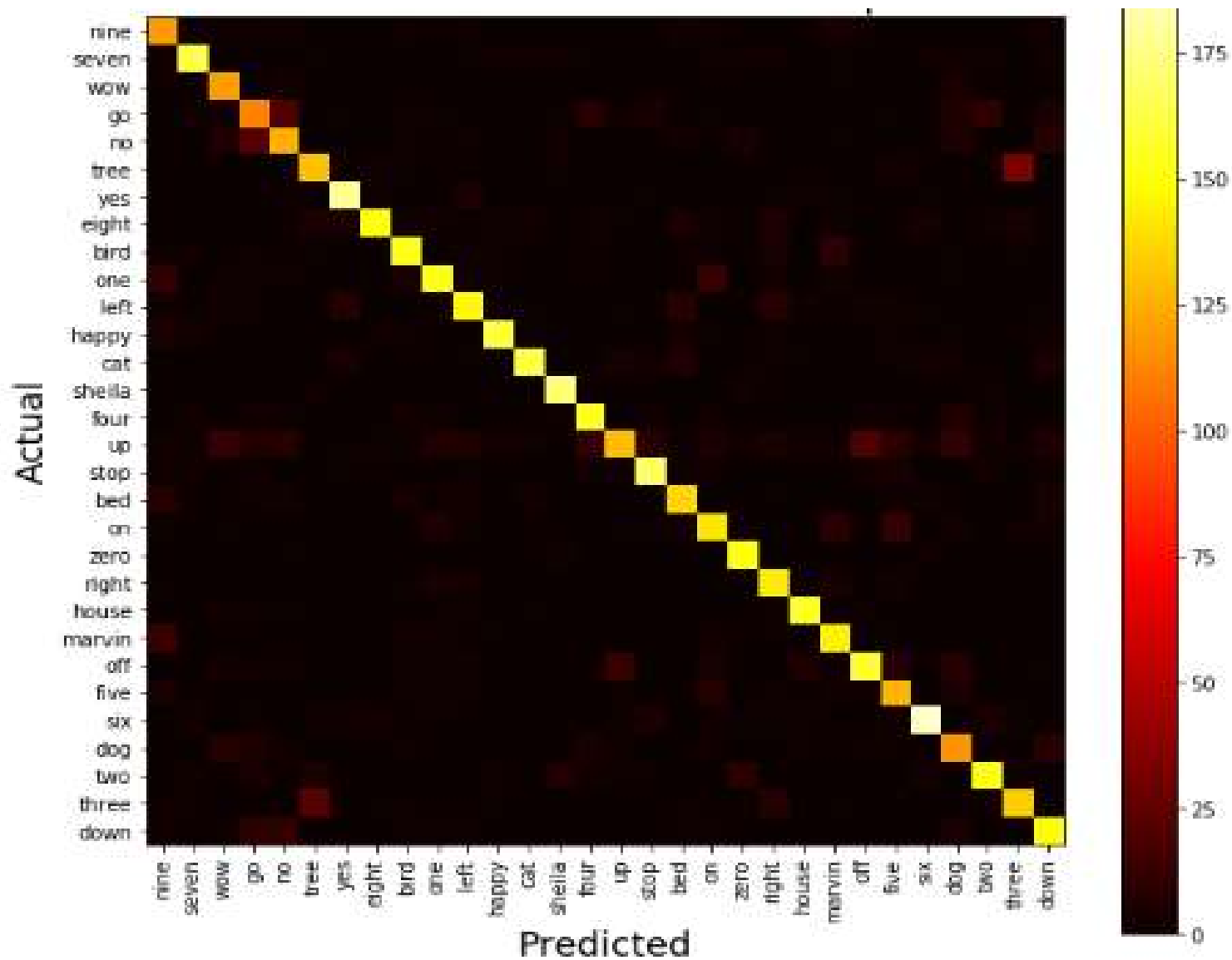




Class	Precision	Recall	F1-score
29 (142)	0.65	0.52	0.76
28 (159)	0.69	0.77	0.73
27 (173)	0.78	0.82	0.80
26 (146)	0.62	0.75	0.67
25 (258)	0.93	0.76	0.83
24 (144)	0.64	0.81	0.72
23 (209)	0.86	0.78	0.82
22 (185)	0.77	0.78	0.77
21 (161)	0.90	0.96	0.93
20 (177)	0.81	0.86	0.84
19 (152)	0.87	0.99	0.93
18 (183)	0.75	0.79	0.77
17 (195)	0.82	0.73	0.77
16 (192)	0.85	0.86	0.86
15 (128)	0.89	0.43	0.58
14 (171)	0.80	0.89	0.85
13 (188)	0.93	0.93	0.93
12 (195)	0.89	0.81	0.85
11 (171)	0.89	0.89	0.89
10 (227)	0.85	0.68	0.76
9 (182)	0.76	0.77	0.77
8 (165)	0.82	0.92	0.87
7 (215)	0.96	0.72	0.82
6 (171)	0.81	0.95	0.88
5 (171)	0.74	0.73	0.74
4 (182)	0.72	0.73	0.73
3 (153)	0.58	0.67	0.62
2 (127)	0.70	0.91	0.79
1 (179)	0.88	0.90	0.89
0 (149)	0.67	0.82	0.74

Deep CNN 2D

- Training and Classification Report for Deep CNN with 20+ layers
 - Layers include
 - 2D Convolution
 - Batch Normalization
 - Dropout
 - Dense Layers
 - Training Accuracy
 - 96%
 - Model score 81%



Conclusions

- Supervised Learning techniques performed poorly. This is probably due to the feature space chosen in addition to the large variability in the dataset
- 1D Convolutional Neural Networks and networks with fewer layers also showed poor performance
- Deep 2D Convolutional Neural Network with MFCC features performed the best giving 81% accuracy
- The CWT coefficient features for the Deep 2D CNN also provided low accuracy
 - This is probably due to having to restrict the signal by using interpolation to reduce time sample size
- Hyperparameter tuning
 - Hyperparameter tuning was used via the hyperopt library in python
 - After optimization of the hyperparameter space the models were run.
 - Issues with the choice of which parameters to vary as well as memory limitations restricted an exhaustive search of the parameter space