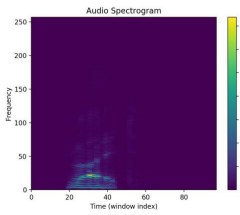
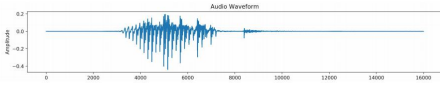


## Objective

Speech recognition is allowing more accessible interactions with agents. The use of deep learning networks has produced high accuracy in speech recognition. We look at using methods from image processing applied to traditional speech recognition systems, such as the use of 2D convolutional and pooling layers for time and frequency invariance in the speech patterns. Our system takes audio waveforms and provides probabilities of 12 different command classes which can be acted on by the agent. Our focus is on efficiency of computation and the latency to determine if an utterance is of a certain class.

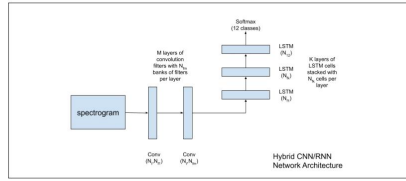
## Dataset

The dataset is taken from the Kaggle TensorFlow Speech Command Recognition Challenge. This consists of 30 classes of spoken commands of which 10 commands are to be recognized. The input data is approximately 1 second of audio sampled at 16KHz. The dataset is 64,727 examples and split into training/validation/test set. The examples are fairly uniformly distributed across the 30 classes. Input data is normalized and converted to spectrogram data to be used as input to our learning network. The FFT and time window are 512 samples in length with 160 sample stride between windows.

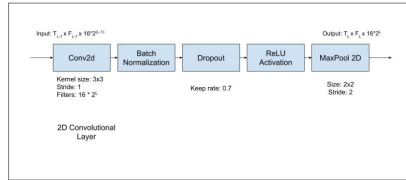


## Architecture

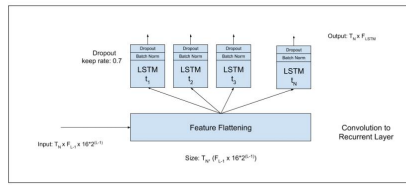
Our model consists of 3 layers of 2D convolutional filtering followed by a layer of feature flattening over each time sample in the sequence. These features are then fed into 3 more layers of LSTM cells. The output is connected through a dense layer to a 12-class softmax layer to determine the probability of spoken command.



Each convolutional layer consists of a 2D convolutional step and a max pooling step for time and frequency invariance.



All of the filter bank outputs and frequency features associated with a single time sample are flattened into a single feature vector and used as the input to the LSTM layers.



## Results

We were able to achieve high accuracy with the 2D convolutional and pooling layers of approximately 95% accuracy. This performed better than our baseline convolutional only and recurrent only networks with 1D convolutional layers.

Architecture	Size	Batch	# Layers	Conv	# LSTM	Other	Parameters	Train	Validation	Accuracy
Baseline	3	1	4	-	-	Batch Norm, Dropout	-	-	-	0.866
Conv1D LSTM	3	1	4	-	-	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	4	-	-	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	1	128	2	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	1	128	2	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv2D LSTM	3	1	3	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	0.876
Conv1D LSTM	3	1	2	128	3	Batch Norm, Dropout	1.39K	2.04K	0.885	