

# Sound Classification

Cough, Laugh, Sigh, Sneeze, Sniff, Throat Clear

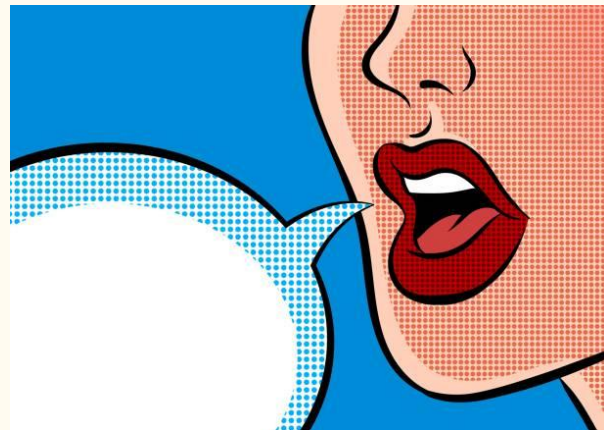
---

**John Guo**

# Motivation

Speech recognition is all around us:

- Siri, Alexa, Google Home
- Youtube or Zoom captioning
- Customer service...
- And many more ([wiki/Speech recognition](https://en.wikipedia.org/wiki/Speech_recognition))



# Motivation

Speech recognition is all around us:

- Siri, Alexa, Google Home
- Youtube or Zoom captioning
- **Customer service...**
- And many more ([wiki/Speech recognition](https://en.wikipedia.org/wiki/Speech_recognition))



# Problem Statement

- Automatic speech recognition (ASR) is hard
  - people make tons of noises unrelated to speech
  - cough, laugh, sigh, sneeze, sniff, throat clear

# Problem Statement

- Automatic speech recognition (ASR) is hard
  - people make tons of noises unrelated to speech
  - cough, laugh, sigh, sneeze, sniff, throat clear
- Goal is to train a neural network that takes in short audio files and classifies it as one of the noises listed above (success: accuracy  $> 90\%$ )

# Problem Statement


- Automatic speech recognition (ASR) is hard
  - people make tons of noises unrelated to speech
  - **cough, laugh, sigh, sneeze, sniff, throat clear**
- Goal is to train a neural network that takes in short audio files and classifies it as one of the noises listed above (success: accuracy  $> 90\%$ )


# Problem Statement

- Automatic speech recognition (ASR) is hard
  - people make tons of noises unrelated to speech
  - **cough, laugh, sigh, sneeze, sniff, throat clear**
- Goal is to train a neural network that takes in short audio files and classifies it as one of the noises listed above (success: accuracy  $> 90\%$ )
- Hope is that this could eventually be paired with ASR

# About our data

- ~21,000 audio files
- ~ 3,000 contributors
- Typically 2-7 seconds long
- Included some info about contributors:
  - Age
  - Gender
  - Country
  - Language

Cough 

Laugh 

Sigh 

Sneeze 

Sniff 

Throat Clear 



I don't think that means what you think it means...


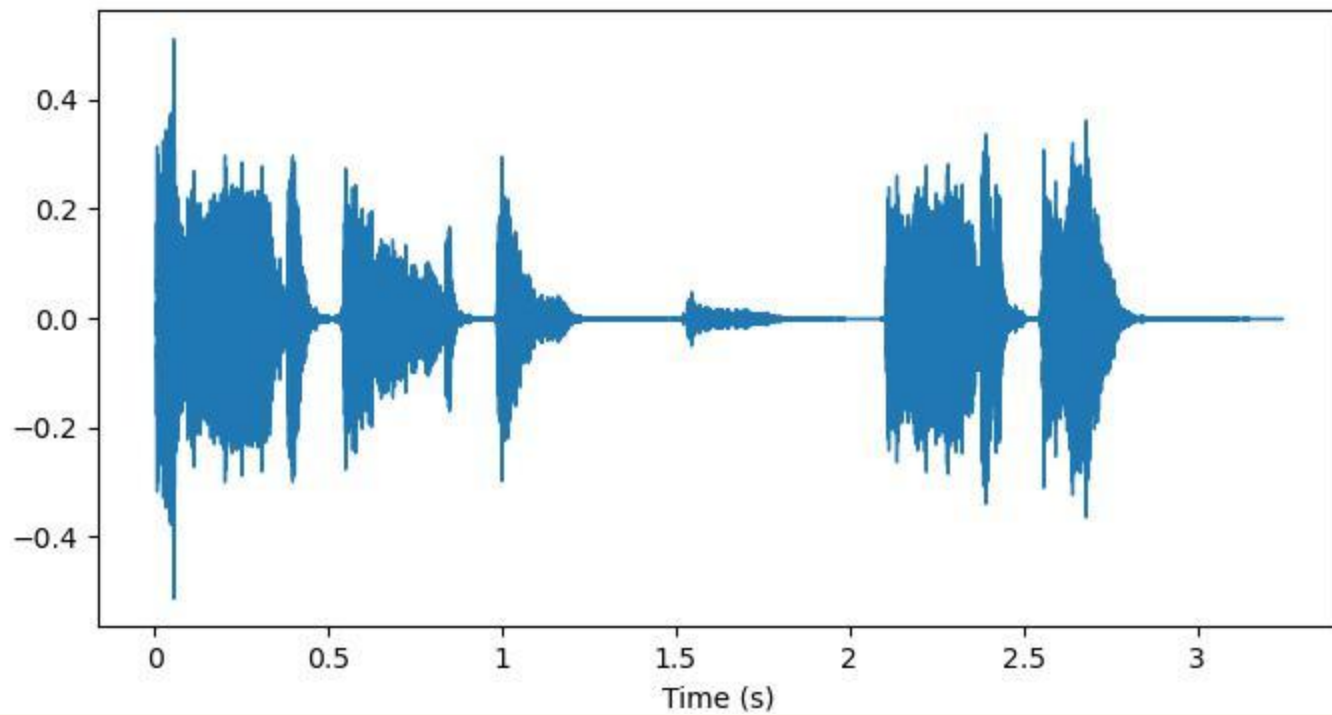
Labeled as a sigh: 

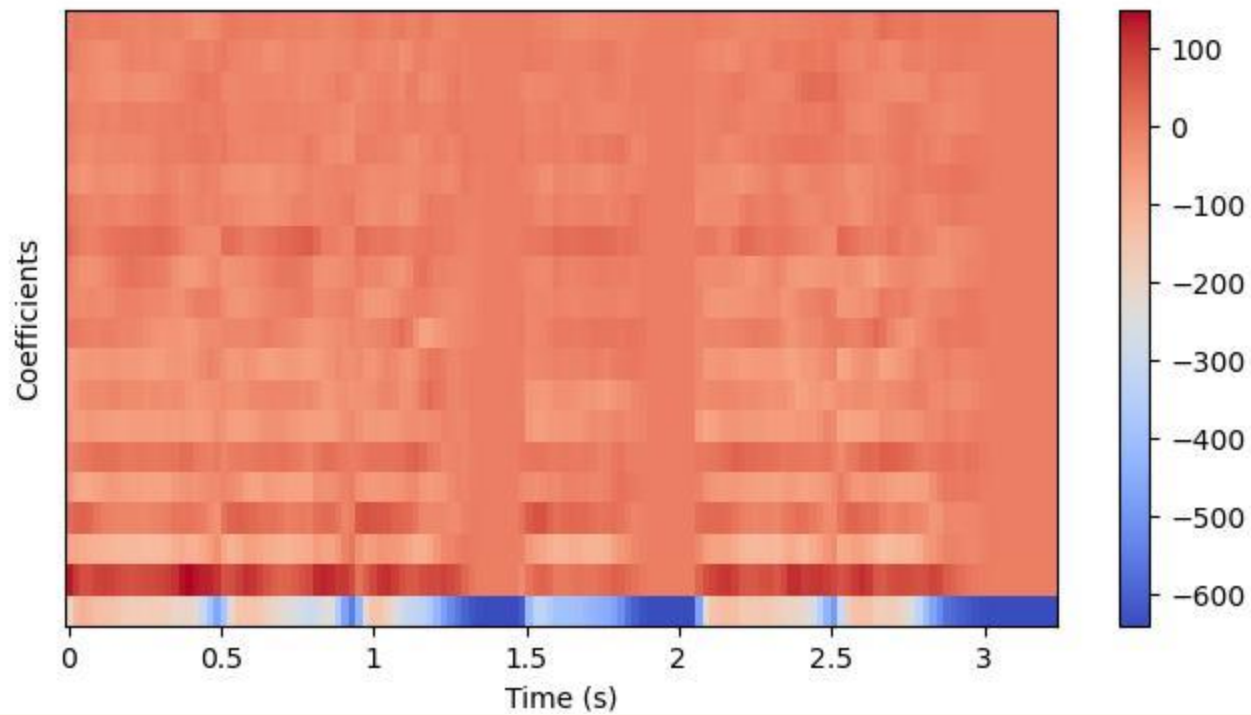


Image provided by Vecteezy

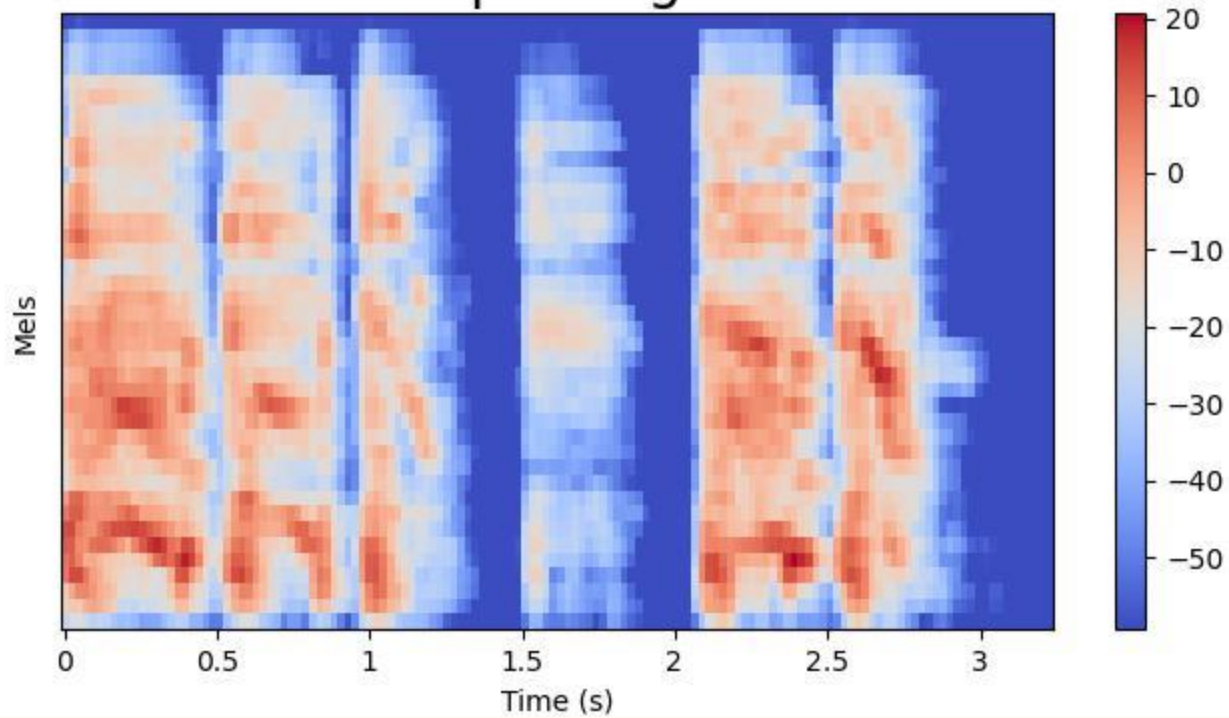
## Waveform



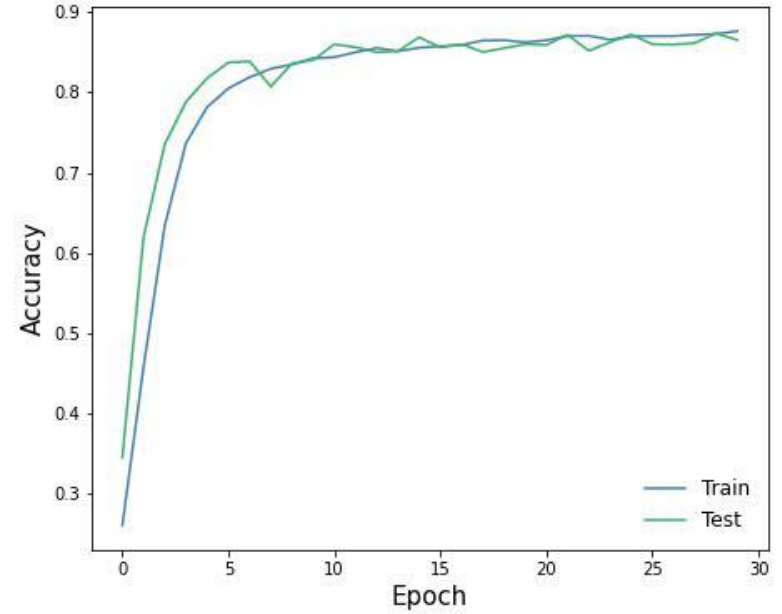
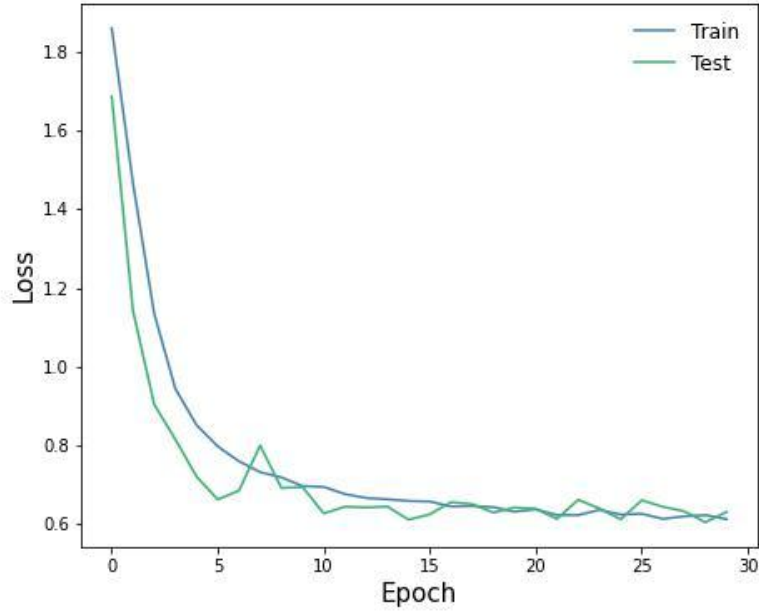
## MFCCs



## Mel Spectrogram

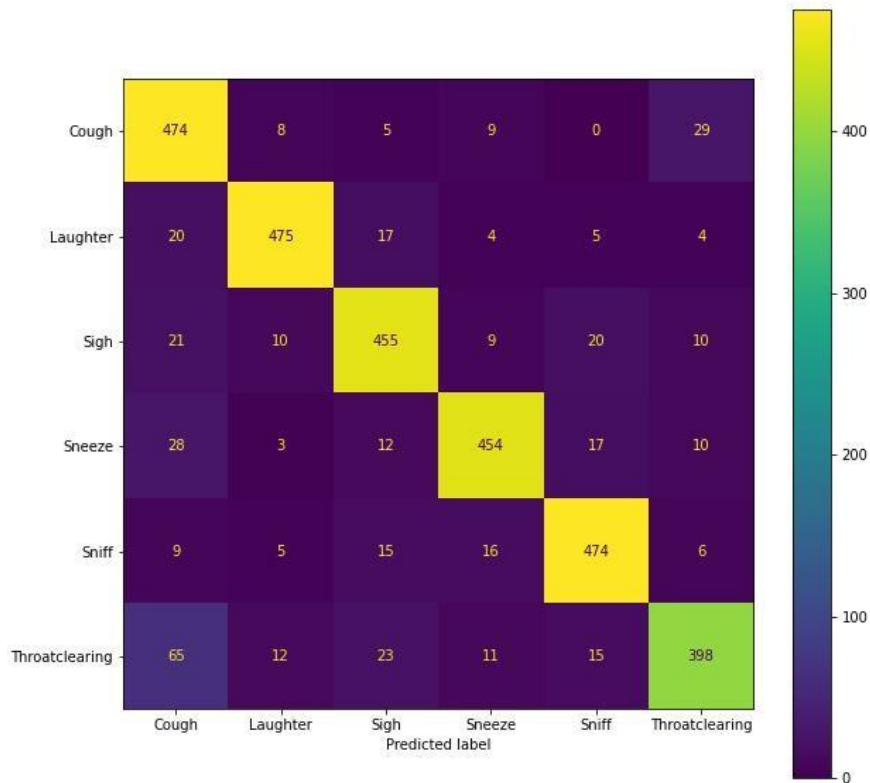


## MFCC CNN Model Performance



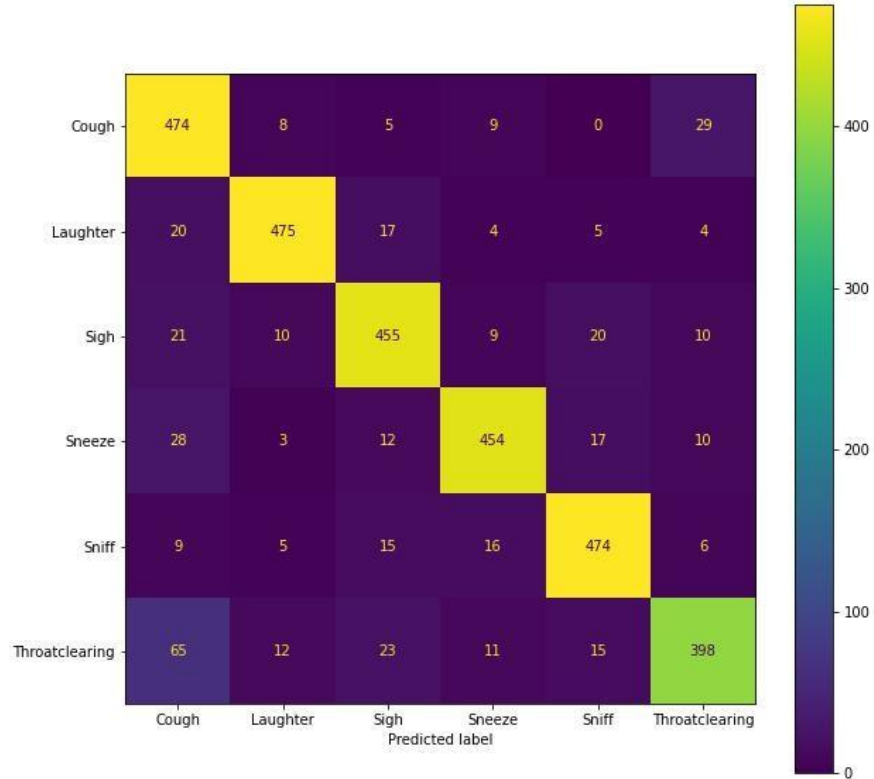
# MFCC Conv. Neural Network Performance

- Big numbers on diagonal



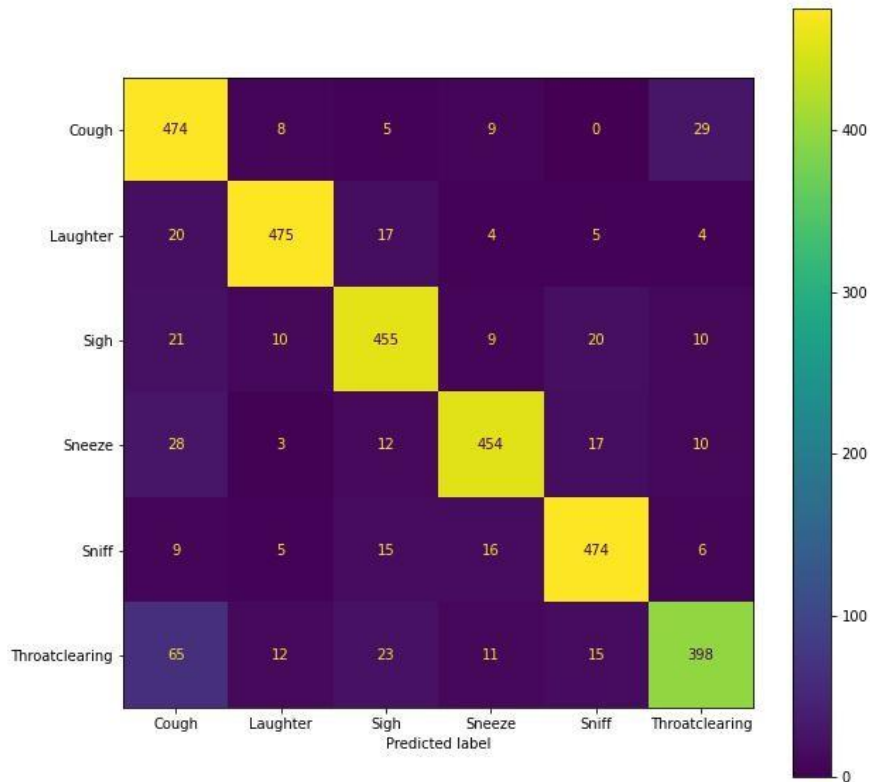
# MFCC Conv. Neural Network Performance

- Big numbers on diagonal
  - Good!



# MFCC Conv. Neural Network Performance

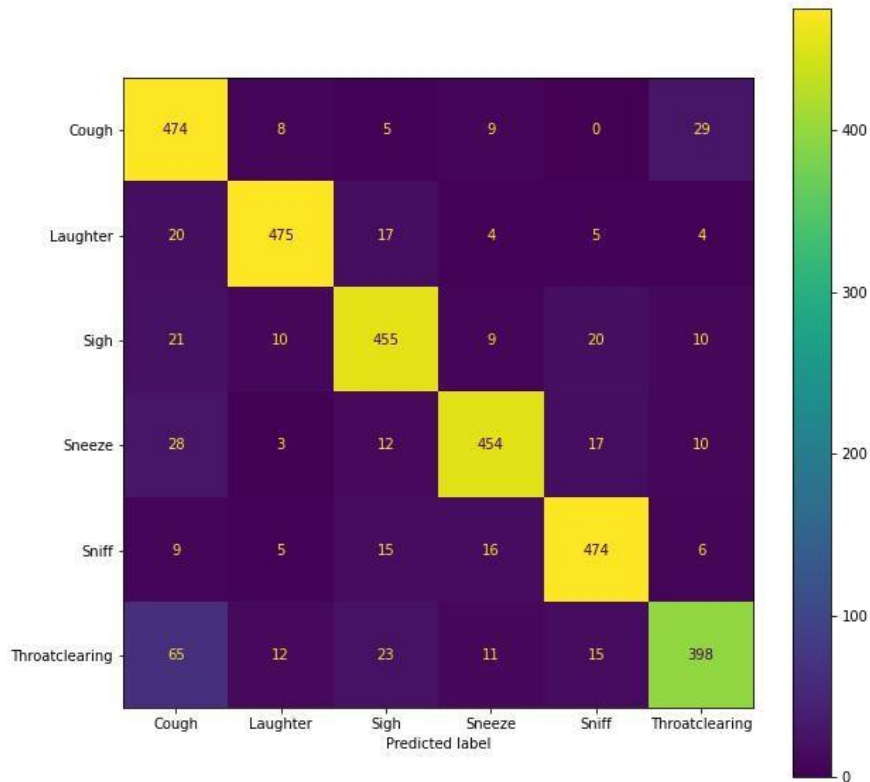
- Big numbers on diagonal
  - Good!
- Small numbers off diagonal





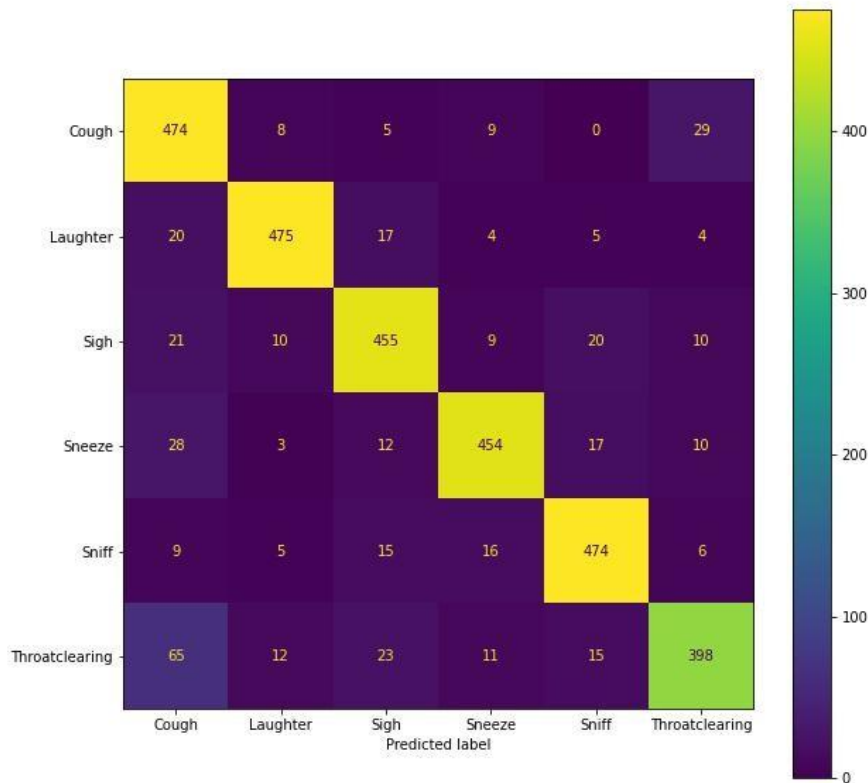
# MFCC Conv. Neural Network Performance

- Big numbers on diagonal
  - Good!
- Small numbers off diagonal
  - Good!

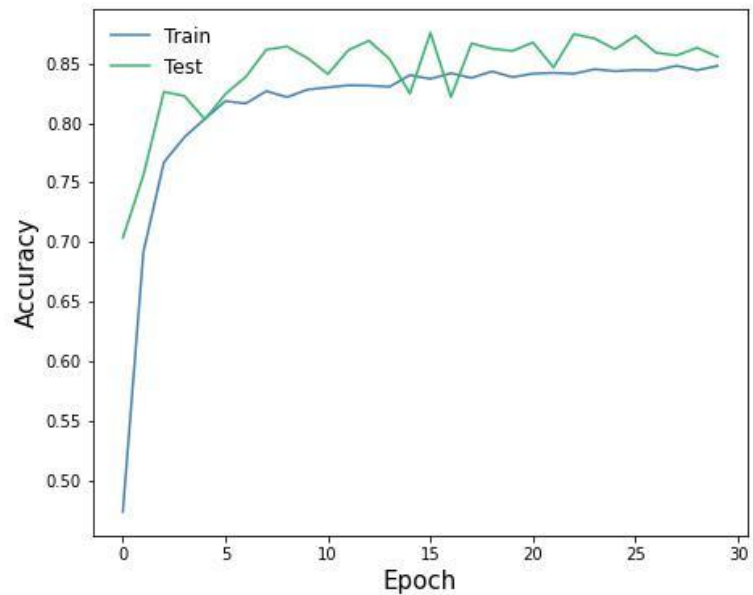
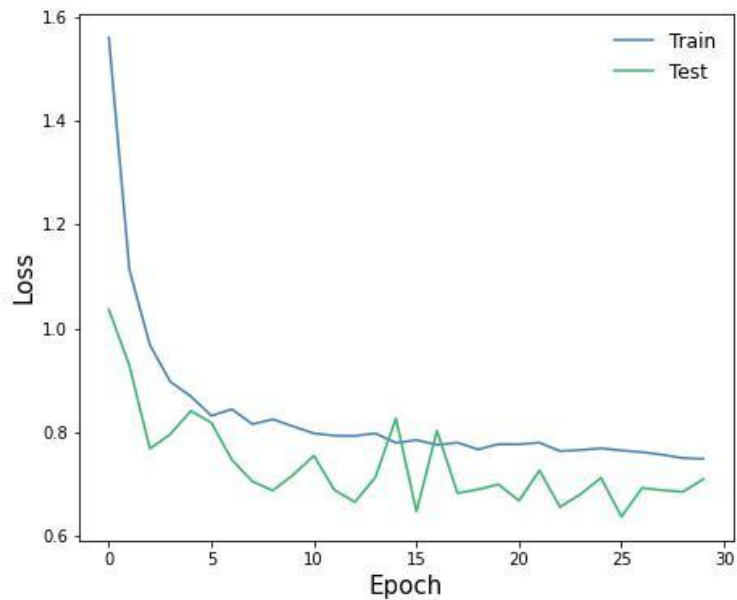


# MFCC Conv. Neural Network Performance

- Big numbers on diagonal
  - Good!
- Small numbers off diagonal
  - Good!
- Labels some throat clears as coughs

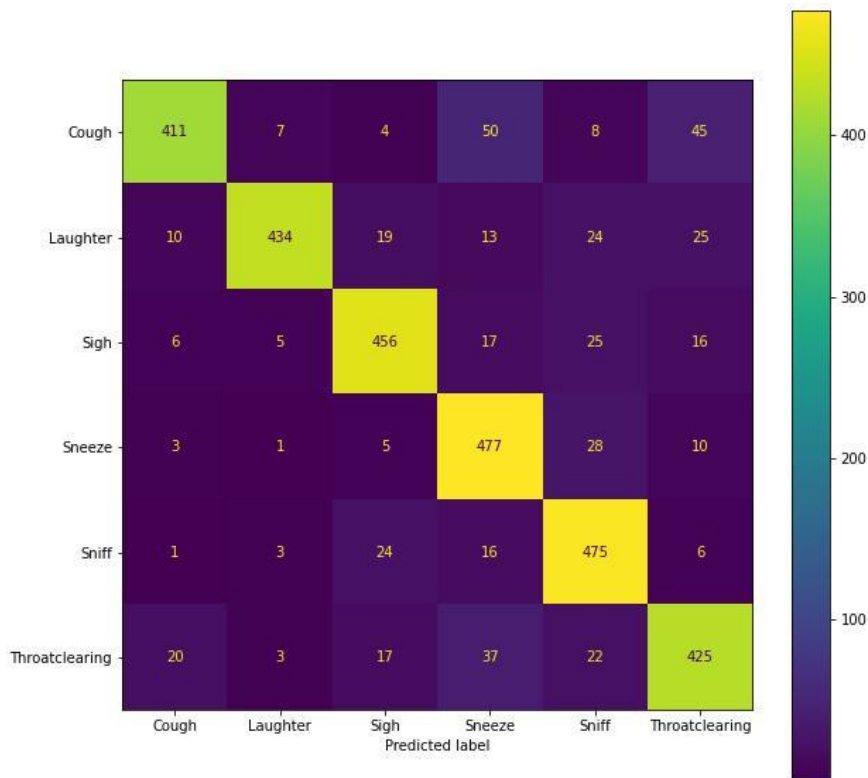


## Mel Spectrogram CNN Model Performance



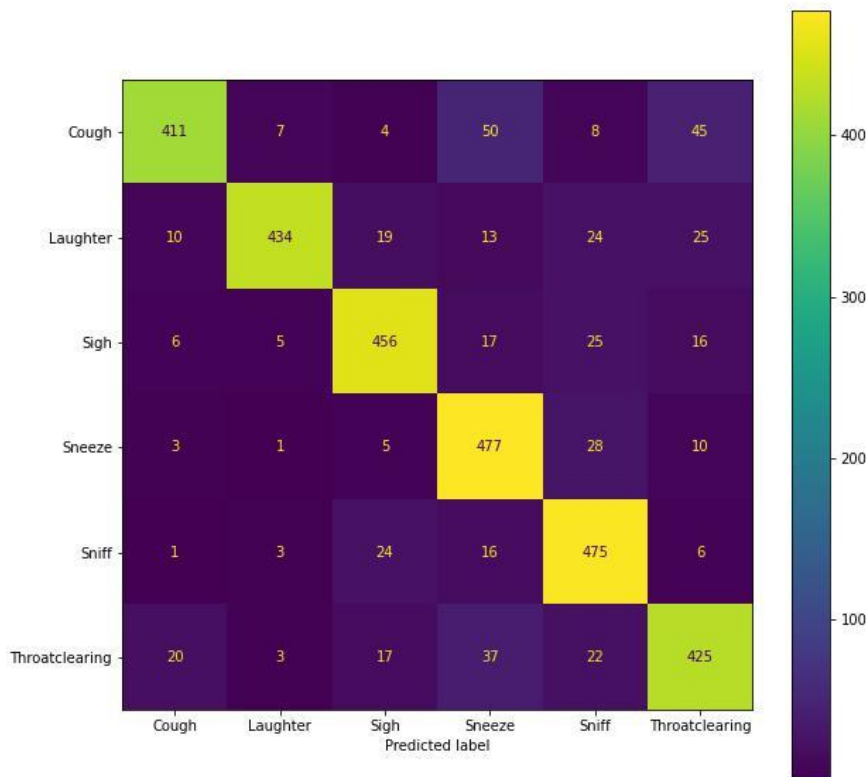
# Mel Spectrogram Conv. Neural Network Performance

- Big numbers on diagonal
  - Good!
- Small numbers off diagonal
  - Good!



# Mel Spectrogram Conv. Neural Network Performance

- Big numbers on diagonal
  - Good!
- Small numbers off diagonal
  - Good!
- Has more general trouble:
  - Sneeze
  - Sniff
  - Throat clearing



Our production model will be...

# Our production model will be...

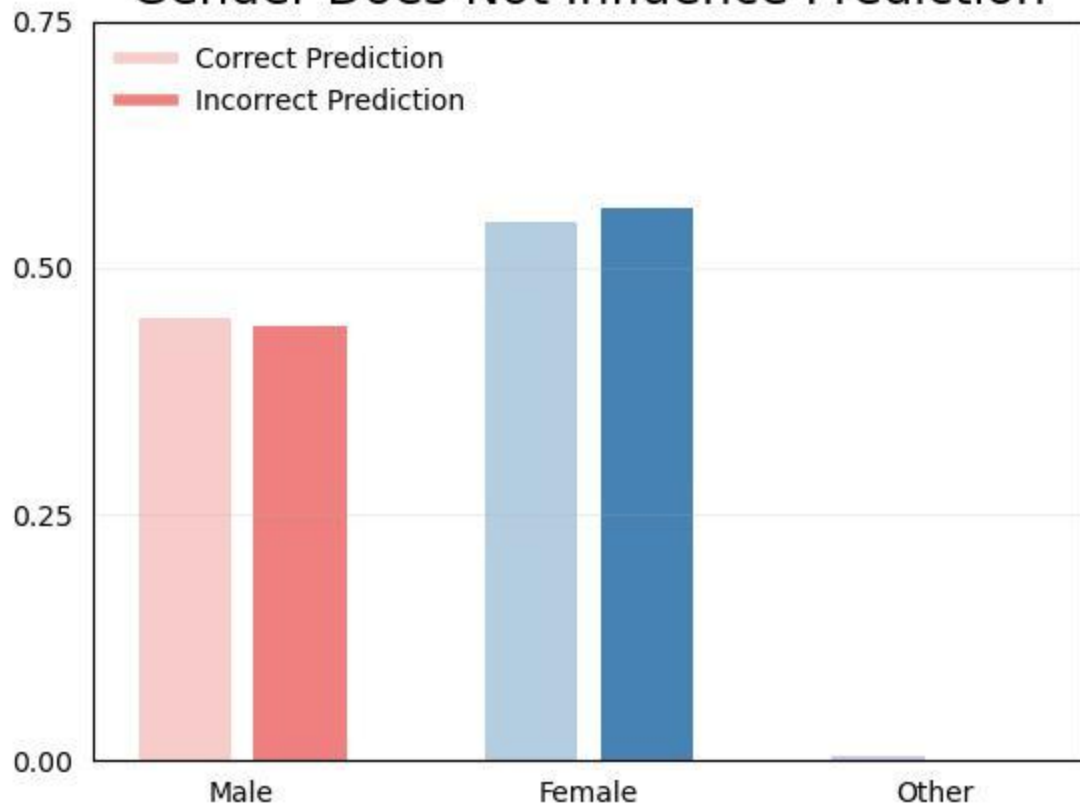
The MFCCs Convolutional Neural Network!

# Our production model will be...

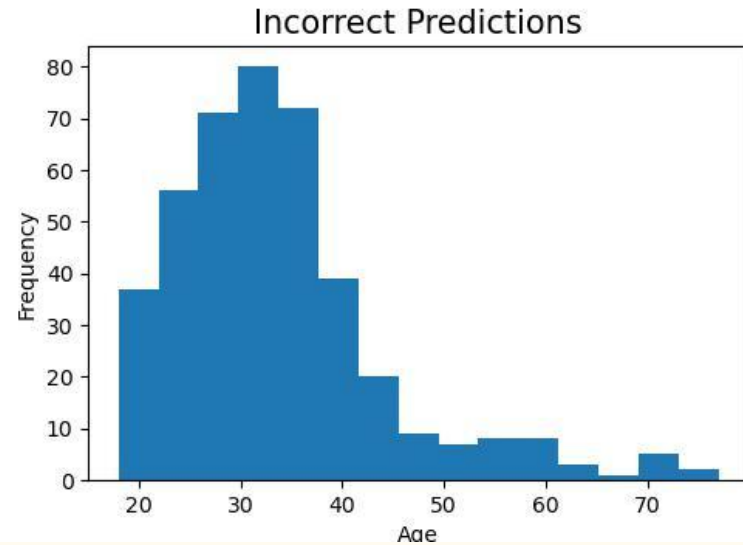
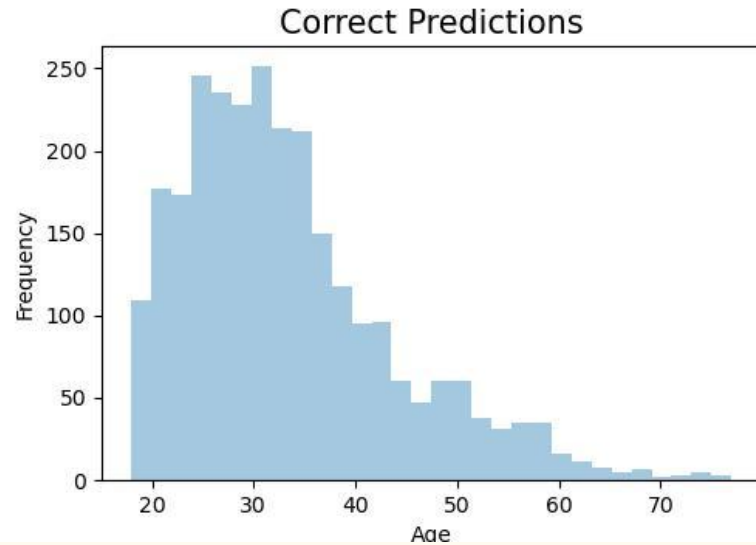
The MFCCs Convolutional Neural Network! (the first one)



## Gender Does Not Influence Prediction

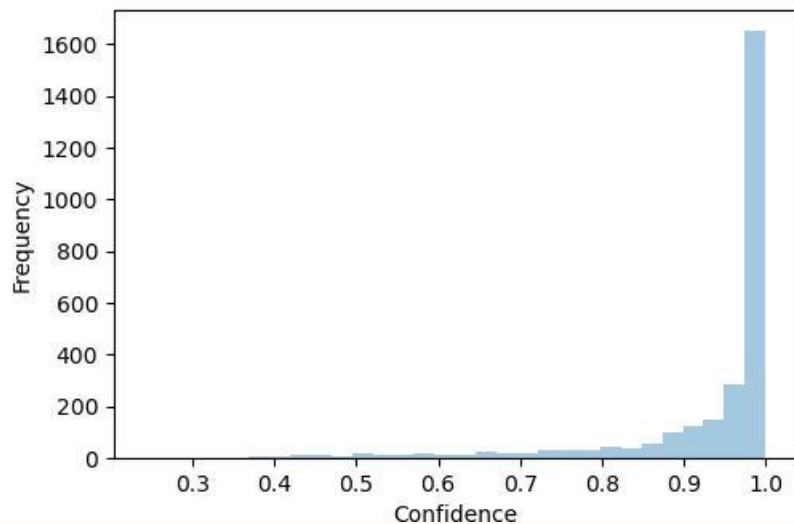


## Distribution of Ages For Correct and Incorrect Predictions

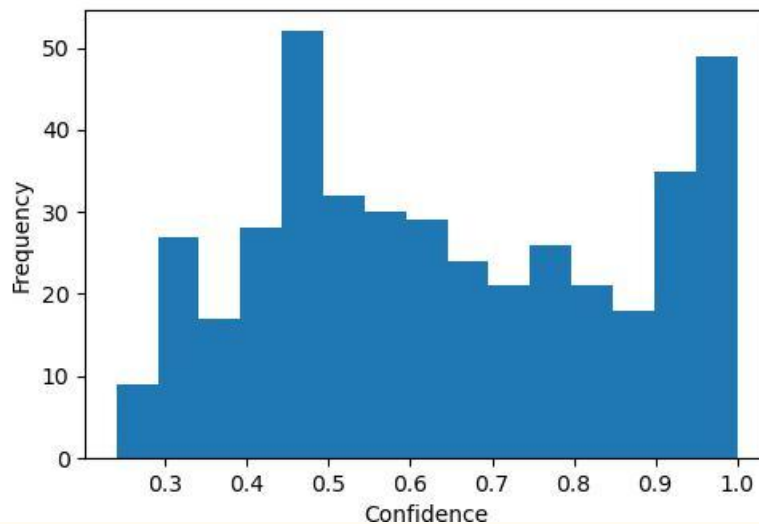


## Distribution of Confidence For Correct and Incorrect Predictions

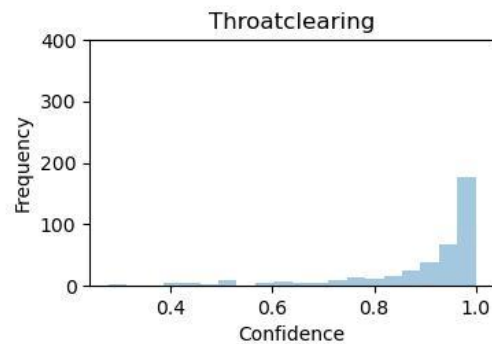
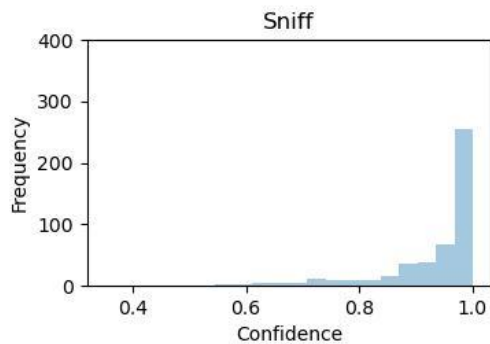
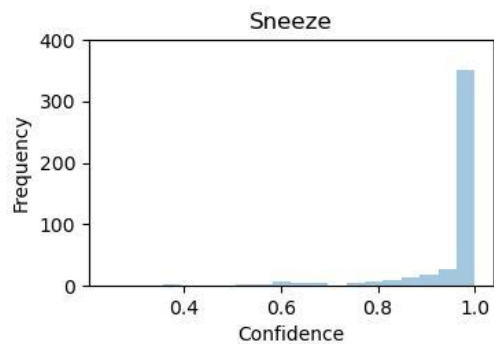
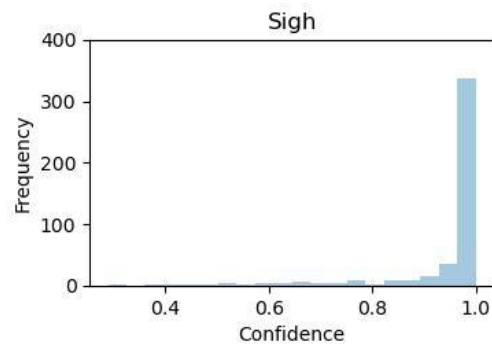
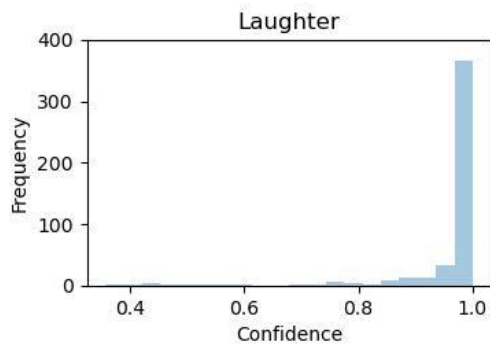
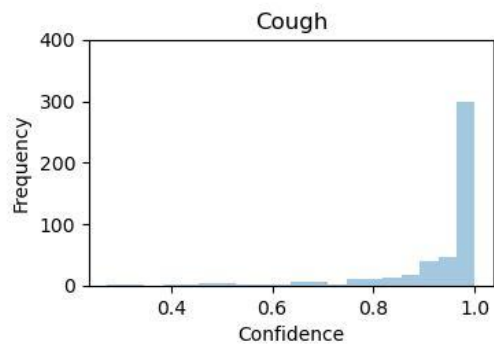
Correct Predictions



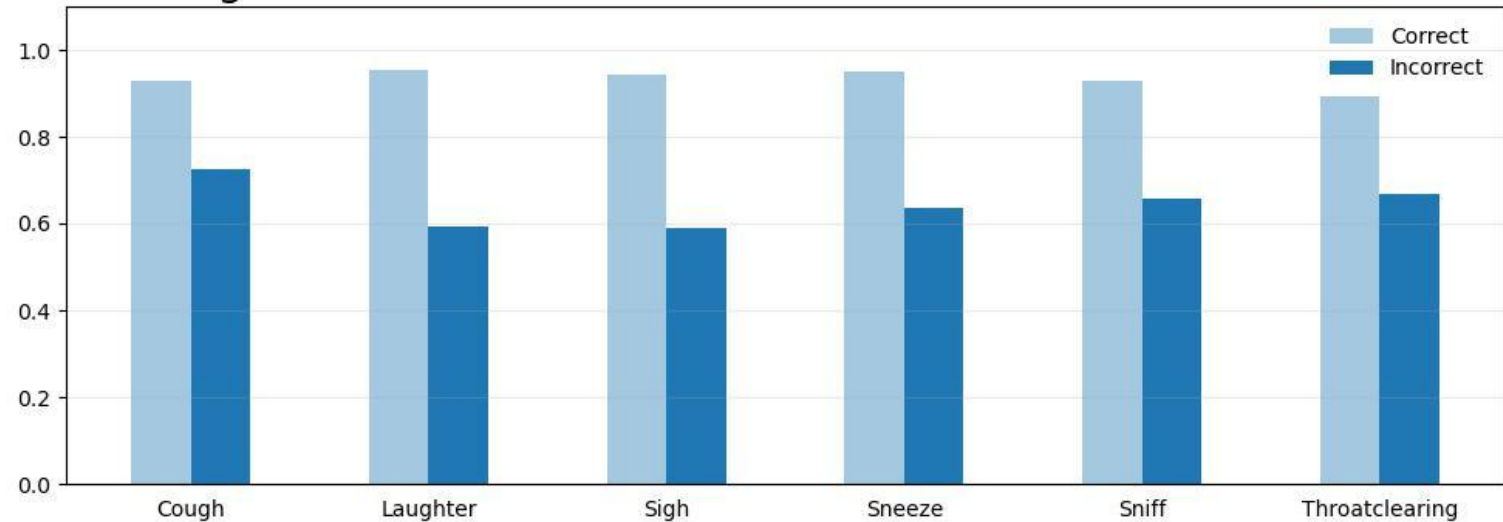
Incorrect Predictions



## Confidence Levels For Correct Predictions



# Average Confidence Levels For Correct and Incorrect Predictions



# Conclusions

- We built a strong model but did not meet our goal of 90% accuracy ( $\sim 87\%$ )
- Developed a Streamlit App to share the model
- Identified where the model had the most trouble (coughs and throat clears)

# Next Steps

- Collect data on other sounds to diversify the classification
- Collect more (properly labeled) audio samples for coughs and throat clears
  - Add this into Streamlit
- Look into applying on longer clips through segmenting data or sliding window
  - Strengthen/support speech recognition

# References & Appreciation

- Valerio Velardo for the education on handling audio data
  - <https://www.youtube.com/@ValerioVelardoTheSoundofAI>
- The dataset:
  - <https://github.com/YuanGongND/vocalsound>
- Kaggle for the free GPU
  - <https://www.kaggle.com>



# To Streamlit!



Image provided by Vecteezy