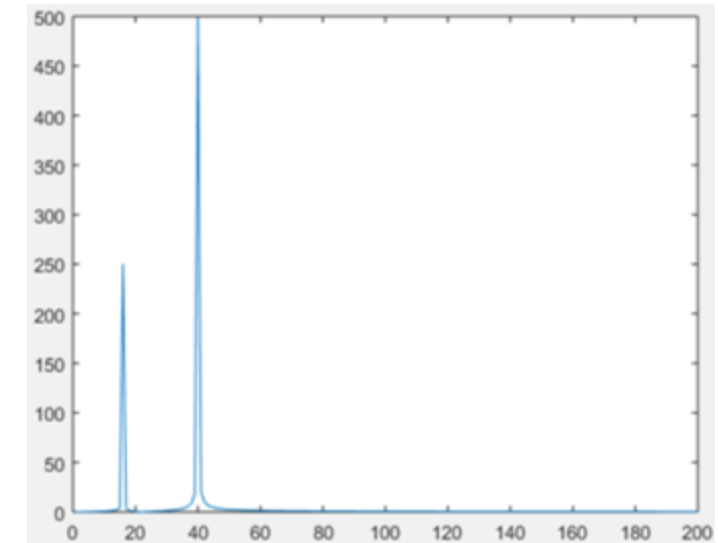| Acquire Raw Data | | Prepare / Clean Data, Visualization | | Feature Engineering | | Pick Model and Hyper-params for Task | | Model Training / Optimization | | Evaluate Model Performance | | Deploy Model |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

- Speech

# Speech

## (Brief Explanation)

# A quick primer on sound

- https://www.youtube.com/watch?v=jveKIYyafaQ
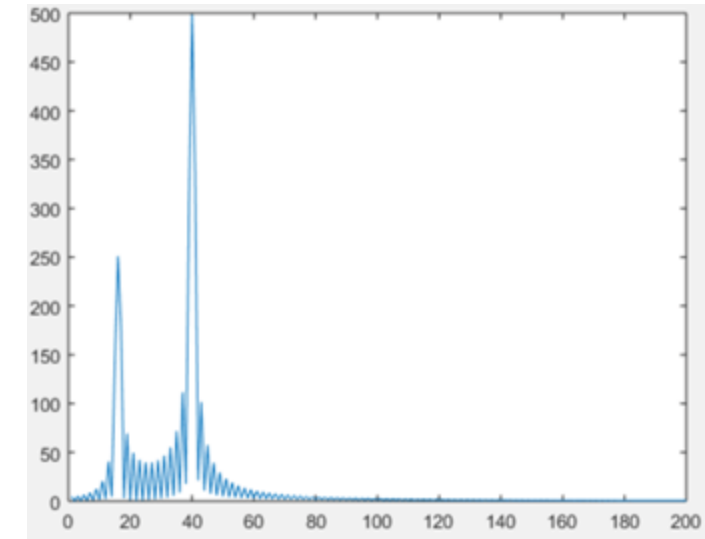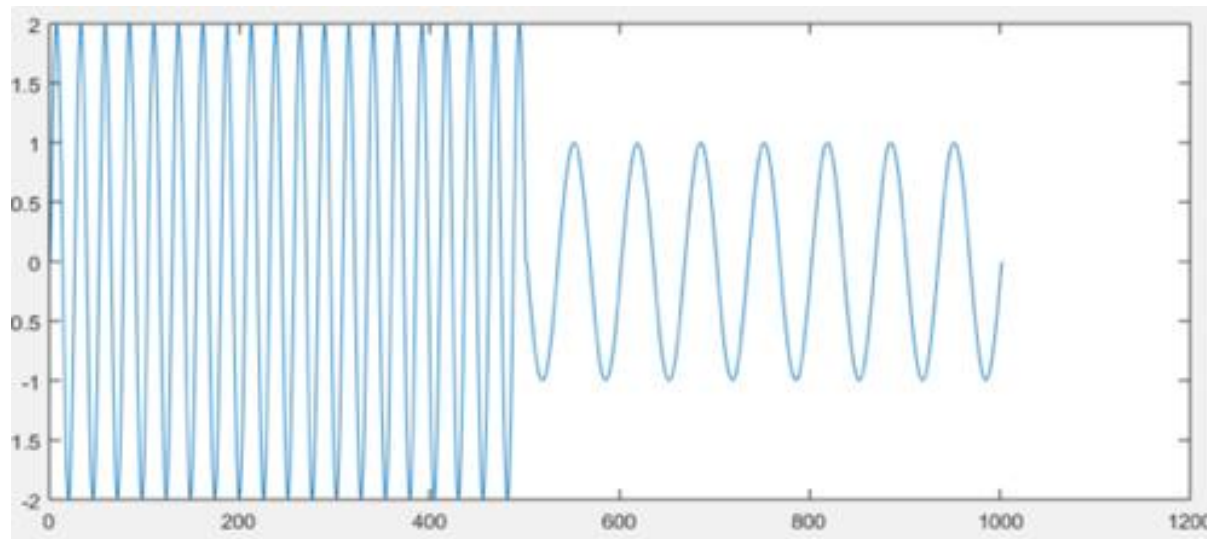
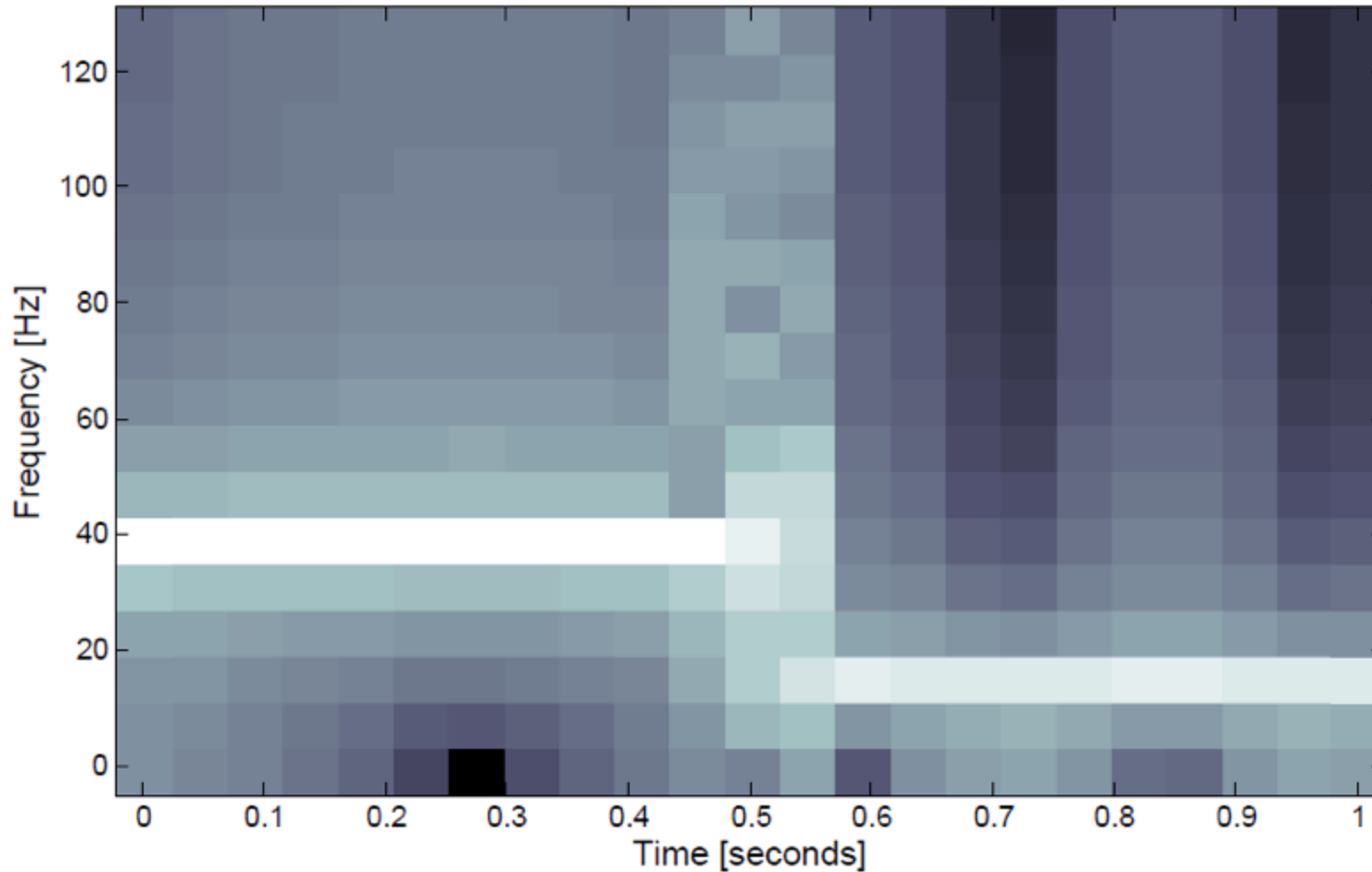# Lets understand Sound Signal

$$f(t) = \sin(2\pi \cdot 39t) + 0.5 \sin(2\pi \cdot 15t)$$

# Example Sound Signal

$$g(t) = \begin{cases} 2 * \sin(2\pi \cdot 39t), & 0 \le t \le 1/2 \\ \sin(2\pi \cdot 15t), & 1/2 < t \le 1 \end{cases}$$
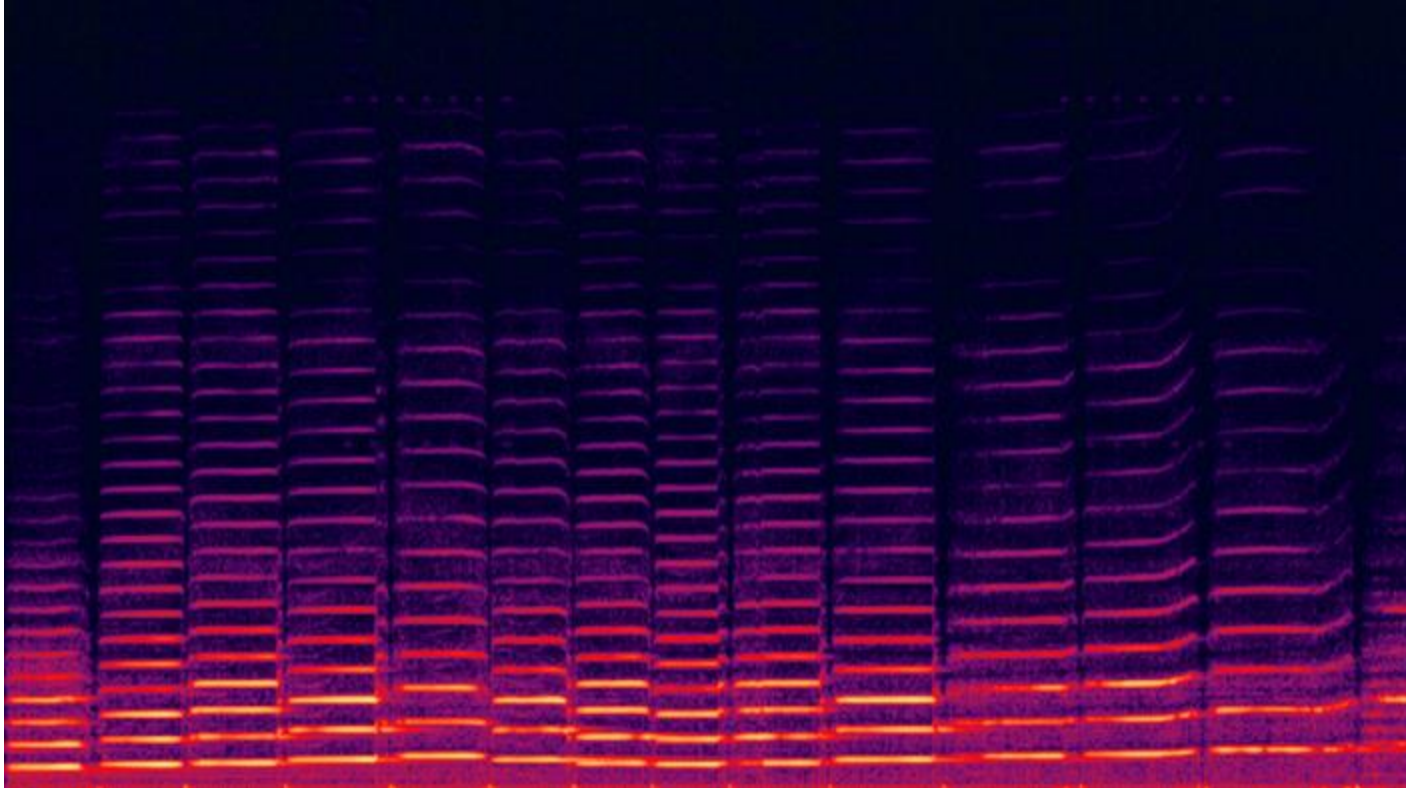
# Spectrogram



Spectrogram of a piecewise monochromatic signal.

Lighter color ▯ greater DFT magnitude

# Spectrogram

# Example Problem

- Sound waves (.wav files)
- 10 short commands ("zero", "one", "two")
- 1 sec duration
- 5000 samples (many people)

# Representations

A: MFCC (Signal processing based; Classical)

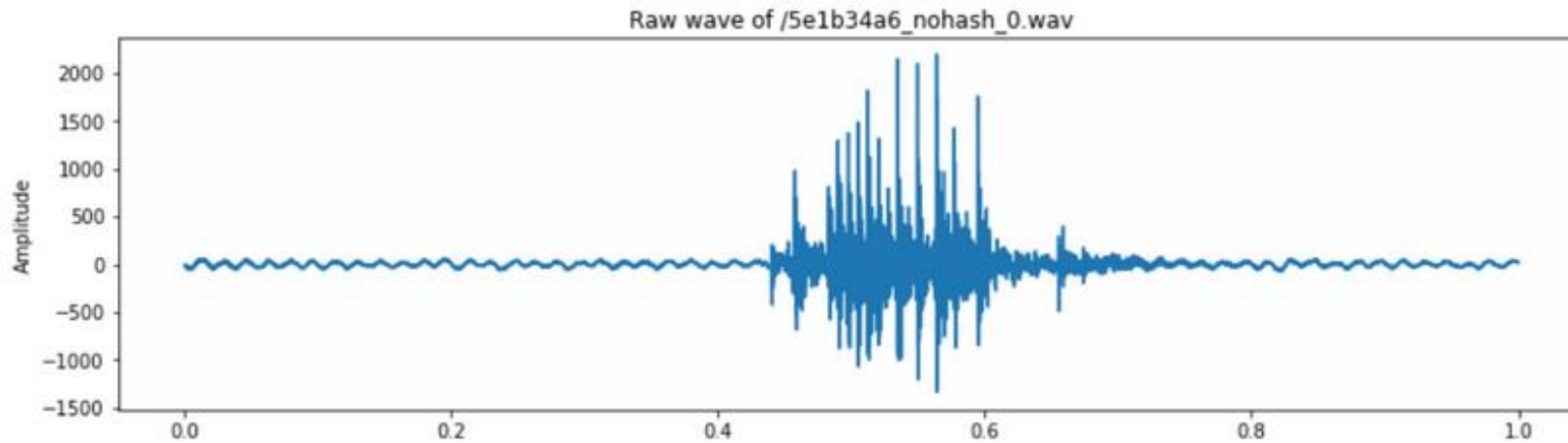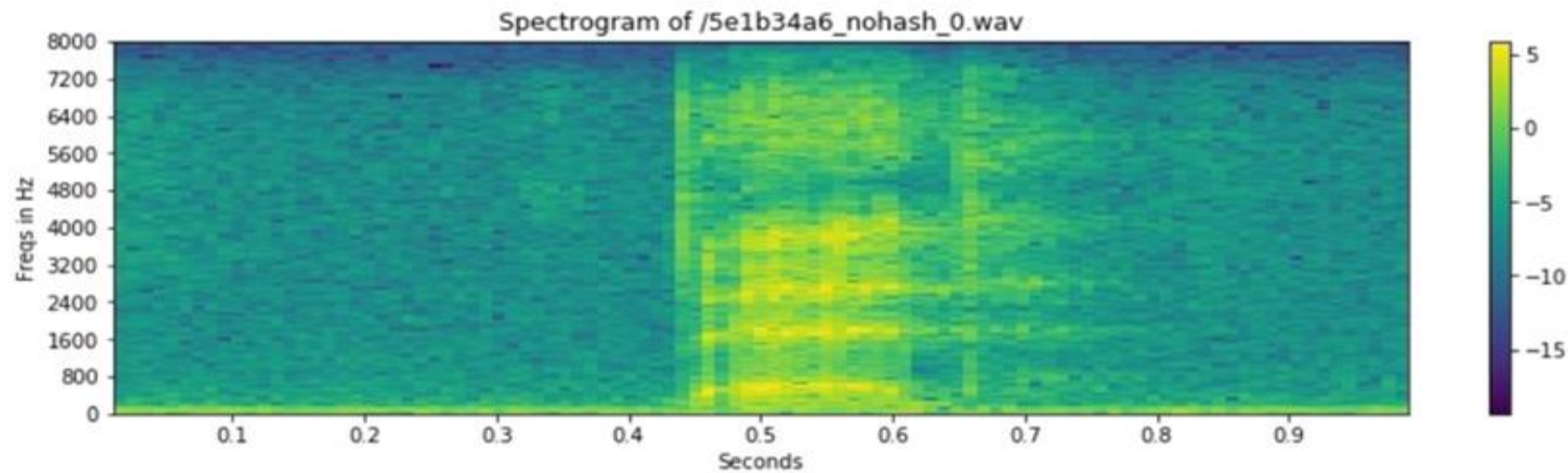- Mel Frequency Cepstral Coefficients

B: CNN Based (Modern)

- VGG Features on the Mel Spectrogram

http://practicalcryptography.com/miscellaneous/machine-learning/guide-mel-frequency-cepstral-coefficients-mfccs/
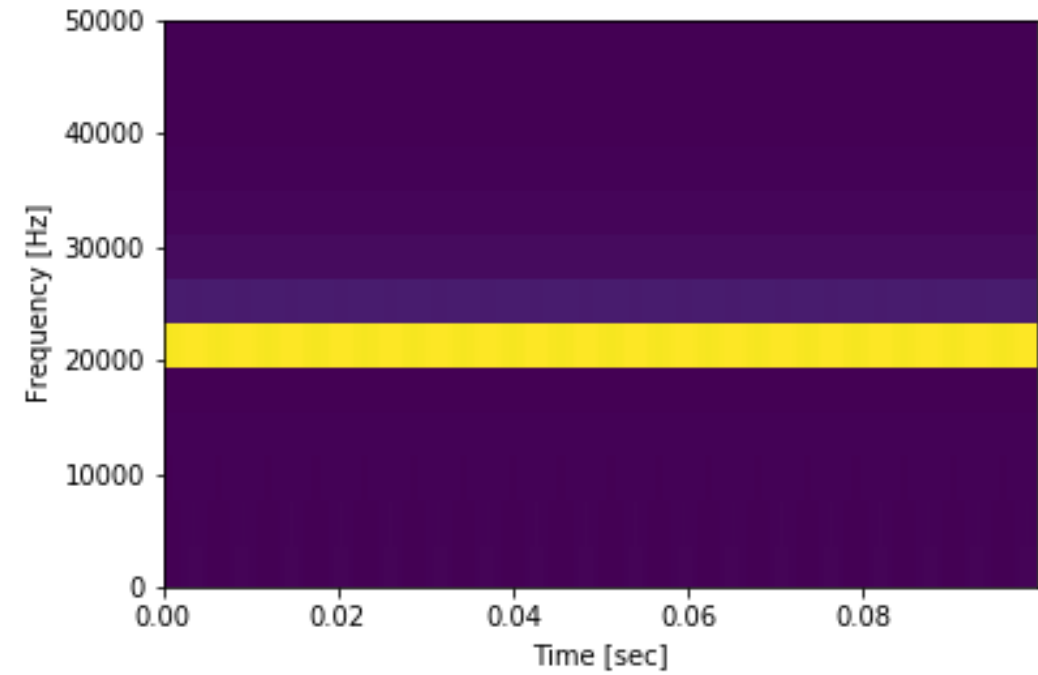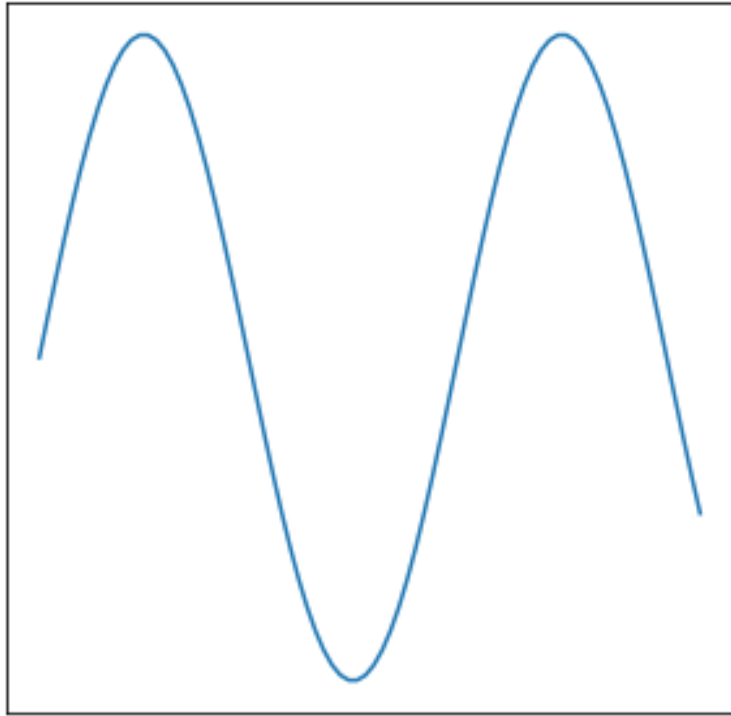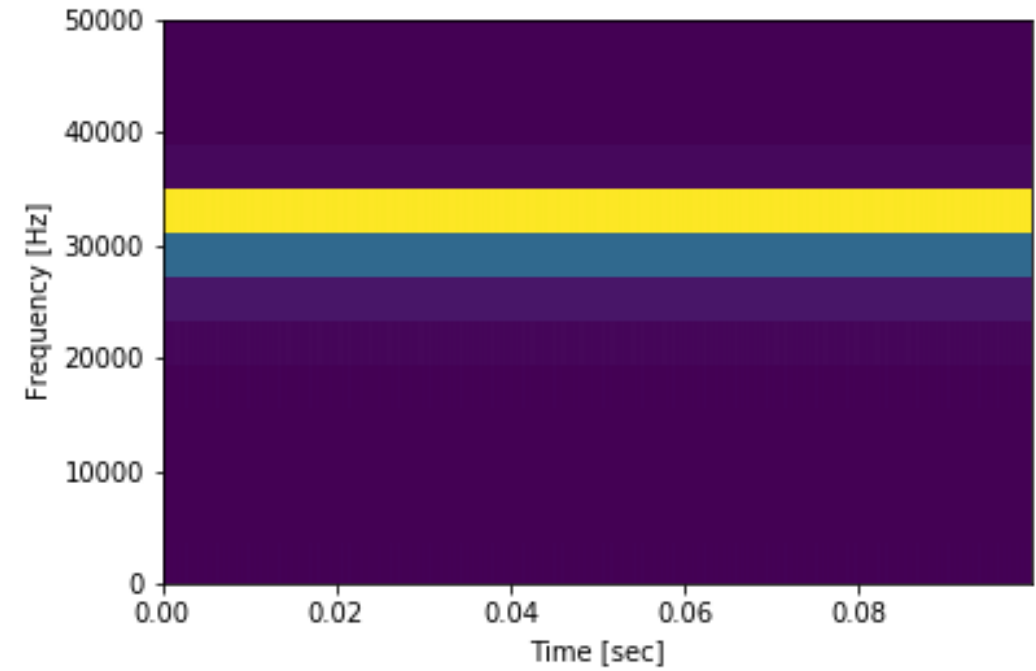
# Classical Feature (MFCC)



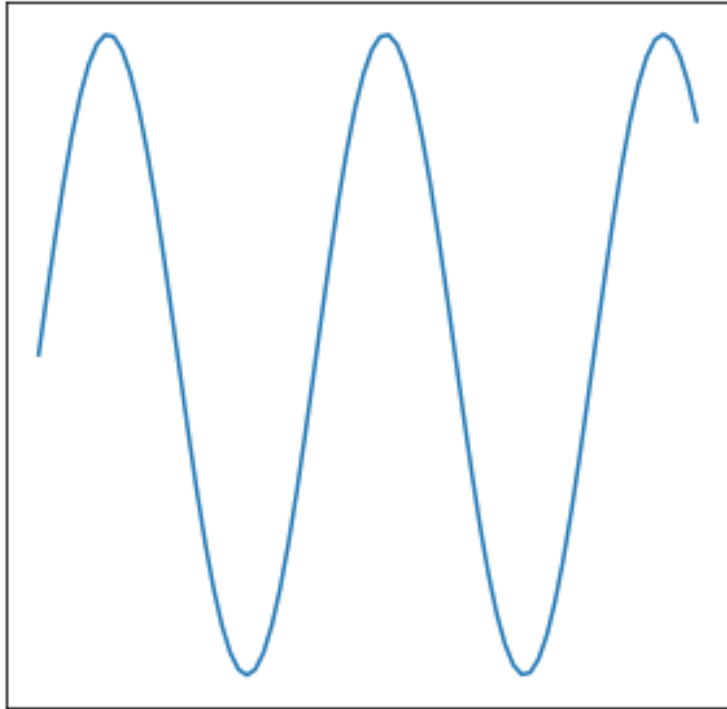Raw wave of /5e1b34a6_nohash_0.wav

**Amplitude Vs Time**



Spectrogram of /5e1b34a6_nohash_0.wav

**Frequency Vs Time**

# Sine wave (f= 20KHz)

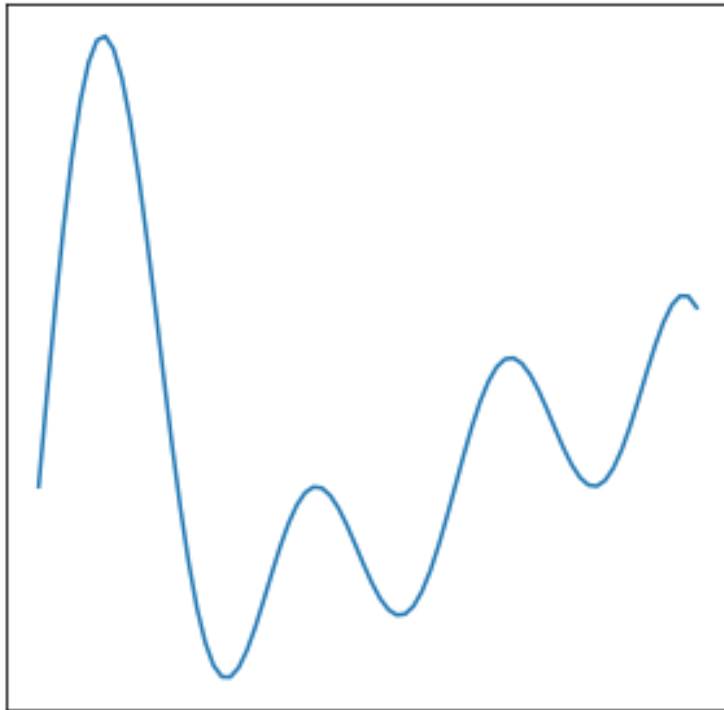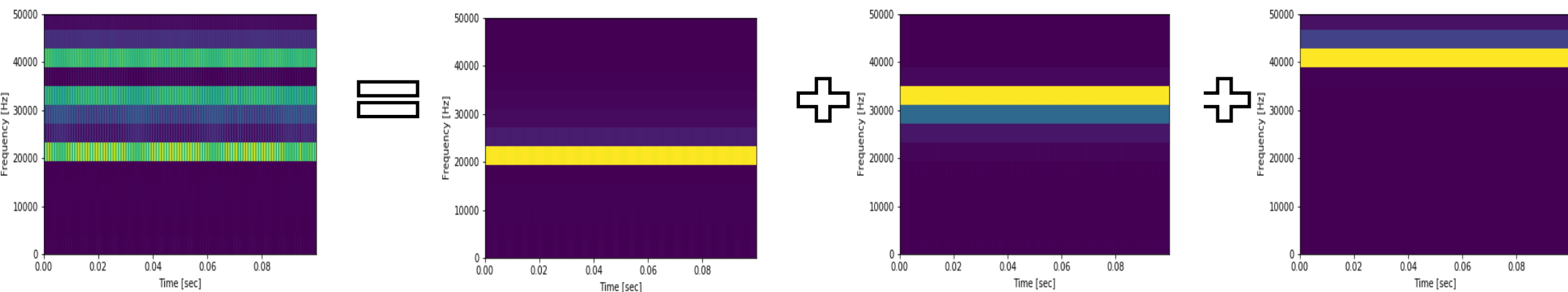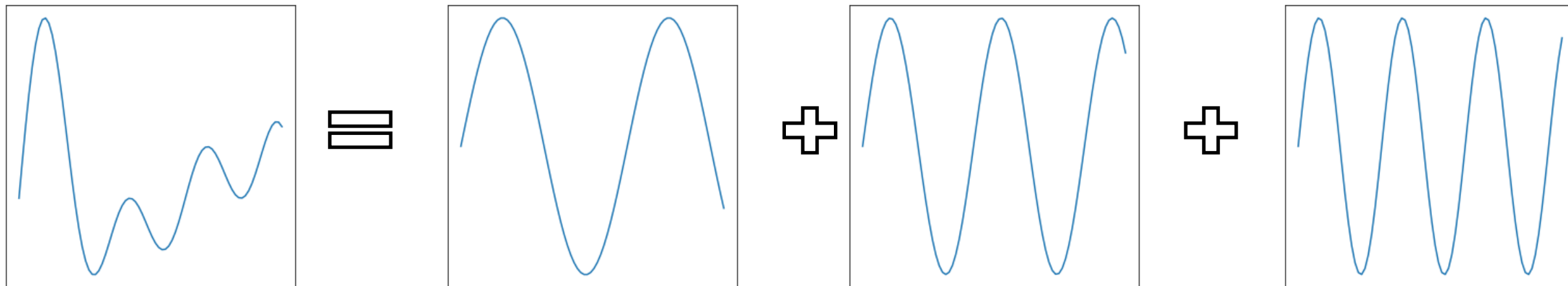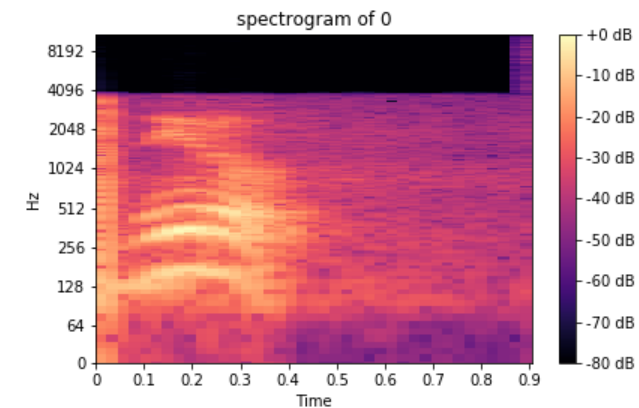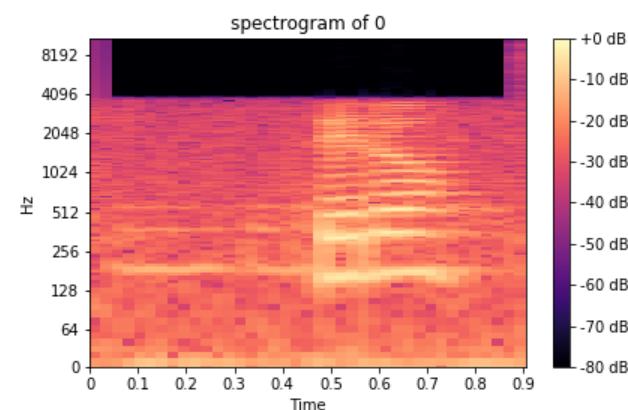# Sine wave (f= 30KHz)

# Sine wave (f= 40KHz)
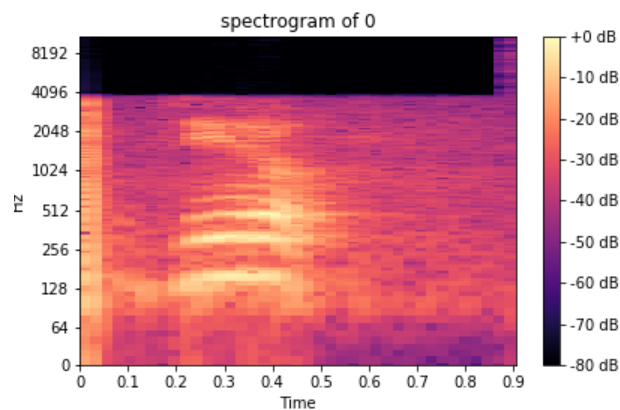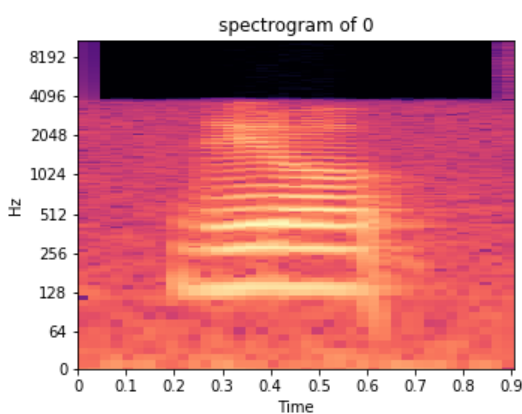
# Any wave is a combination of many sine waves

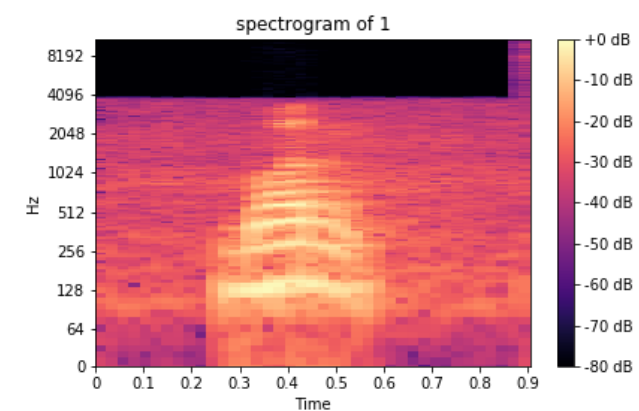# Any wave is a combination of many sine waves

# Example Problem

- Sound waves (.wav files)
- 10 short commands ("zero", "one", "two")
- 1 sec duration
- 5000 samples (many people)

# Utterance of the Word Zero

# Utterance of the Word One

# Utterance of the Word Two



spectrogram of 2

spectrogram of 2

spectrogram of 2

spectrogram of 2

# Features from Mel Spectrogram



Mel spectrogram

**MFCC**
**(Hand coded Classic Features)**

**VGG19-Features**
**(Trained on Mel spectrograms)**

http://practicalcryptography.com/miscellaneous/machine-learning/guide-mel-frequency-cepstral-coefficients-mfccs/

# Example problem



VoxCeleb

*A large scale audio-visual dataset of human speech*

**VoxCeleb2**

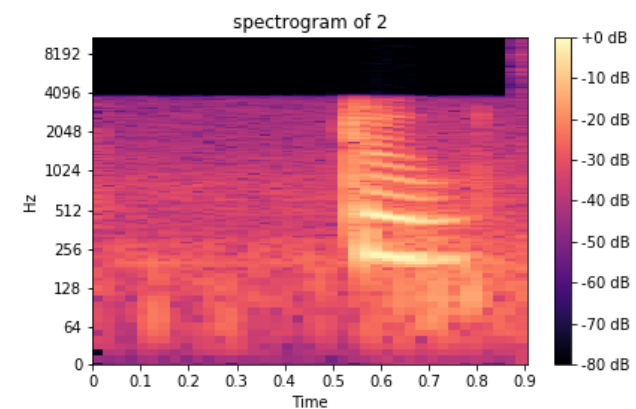*VoxCeleb2 contains over a million utterances for 6,112 identities.*



39%

61%

- Male
- Female



7%

6%

6%

10%

6%

29%

- U.S.A.
- U.K.
- Unknown
- Germany
- India
- France

# Example problem

# Performance on VoxCeleb

| Accuracy | Top-1 (%) | Top-5 (%) |
|---|---|---|
| I-vectors + SVM | 49.0 | 56.6 |
| I-vectors + PLDA + SVM | 60.8 | 75.6 |
| CNN-fc-3s no var. norm. | 63.5 | 80.3 |
| CNN-fc-3s | 72.4 | 87.4 |
| **CNN** | **80.5** | **92.1** |

# Neural Networks + word2vec for text



Sentiment Classification
(Positive / Negative)

# Summary

- Data driven features are now effective for many data.
  - "Learn from some one else's data".
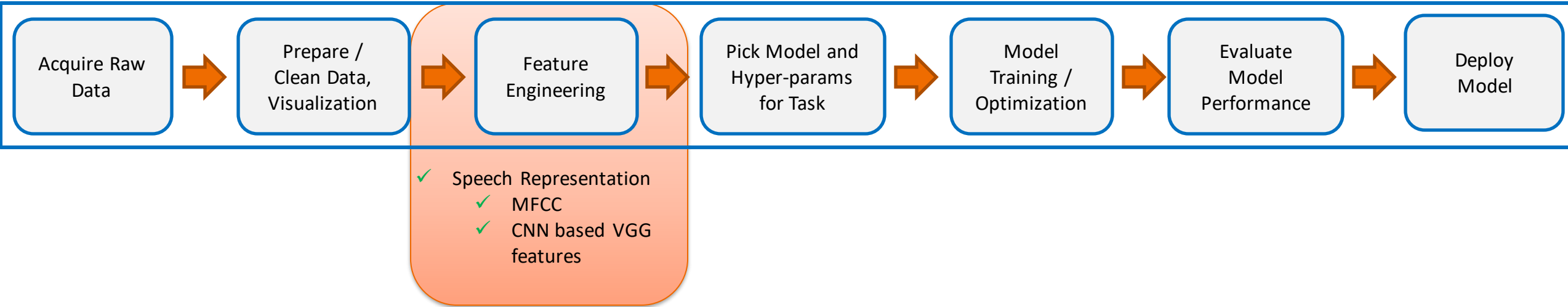  - "Refine to your problem" (more later)
- Many recognition/classification tasks in the image and speech space are reachable.

# Summary



Acquire Raw Data → Prepare / Clean Data, Visualization → Feature Engineering → Pick Model and Hyper-params for Task → Model Training / Optimization → Evaluate Model Performance → Deploy Model

- ✓ Speech Representation
  - ✓ MFCC
  - ✓ CNN based VGG features

# Thanks!!

## Questions?