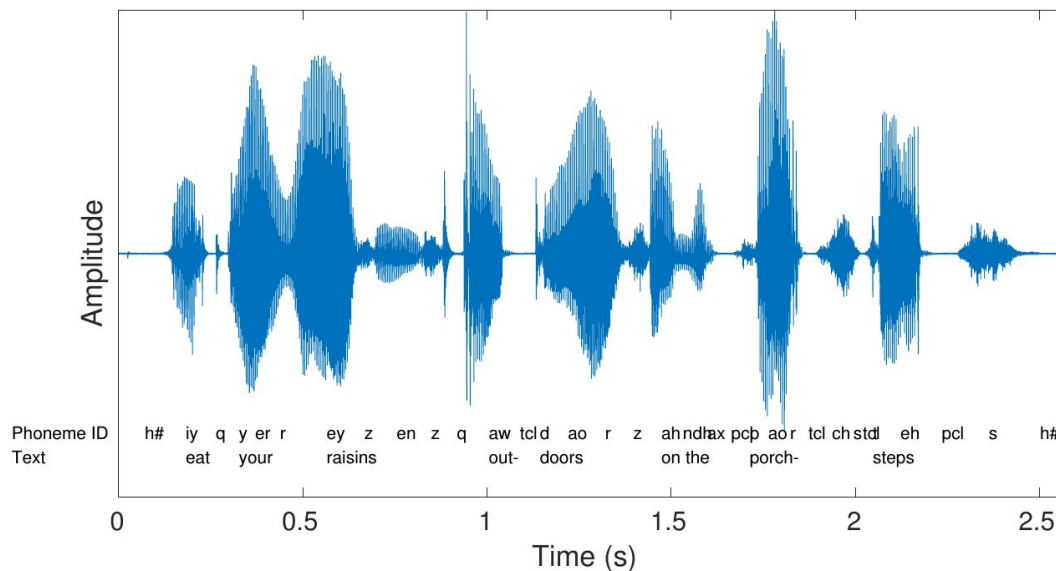


# ASR Features

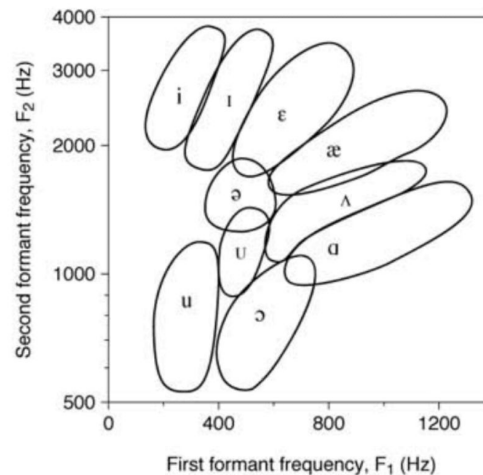
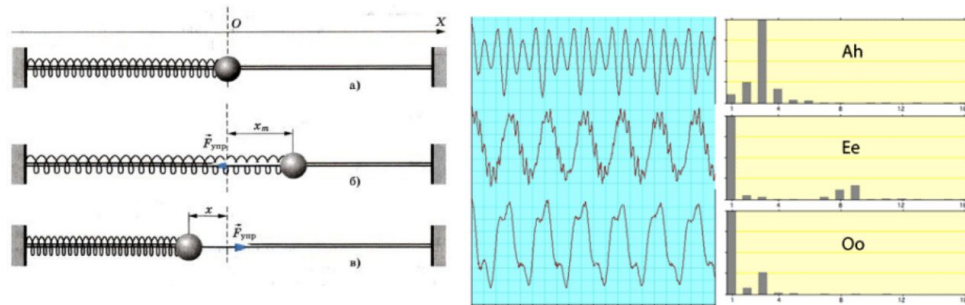
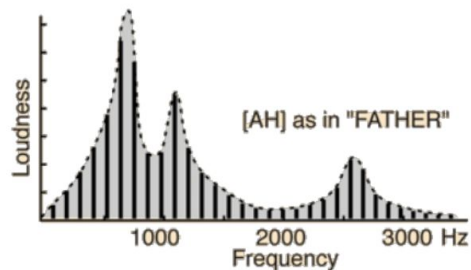
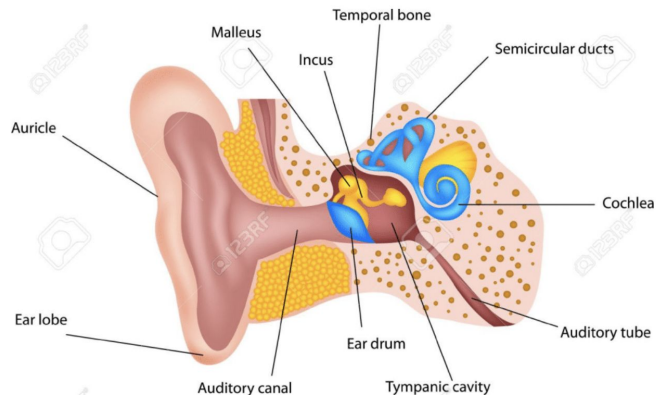
# Task example:

We want to spot some fixed phrase or word.

Input is a simple wav file.



# How do human do it?



# Discrete fourier transform

$$x_n = \frac{1}{N} \sum_{k=0}^{N-1} X_k \cdot e^{i2\pi kn/N} \quad (\text{Eq.3})$$

## Definition [\[ edit \]](#)

The *discrete Fourier transform* transforms a **sequence of  $N$  complex numbers**  $\{\mathbf{x}_n\} := x_0, x_1, \dots, x_{N-1}$  into another sequence of complex numbers,  $\{\mathbf{X}_k\} := X_0, X_1, \dots, X_{N-1}$ , which is defined by

$$\begin{aligned} X_k &= \sum_{n=0}^{N-1} x_n \cdot e^{-\frac{i2\pi}{N}kn} \\ &= \sum_{n=0}^{N-1} x_n \cdot \left[ \cos\left(\frac{2\pi}{N}kn\right) - i \cdot \sin\left(\frac{2\pi}{N}kn\right) \right], \end{aligned} \quad (\text{Eq.1})$$

# Framing

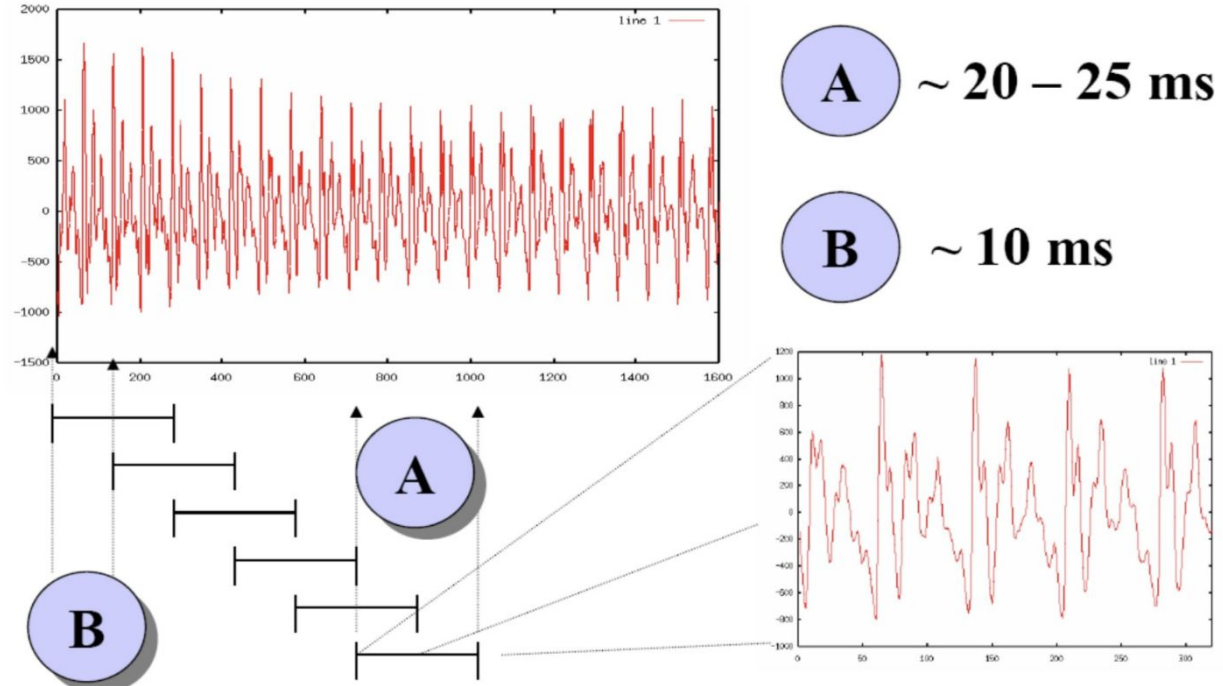
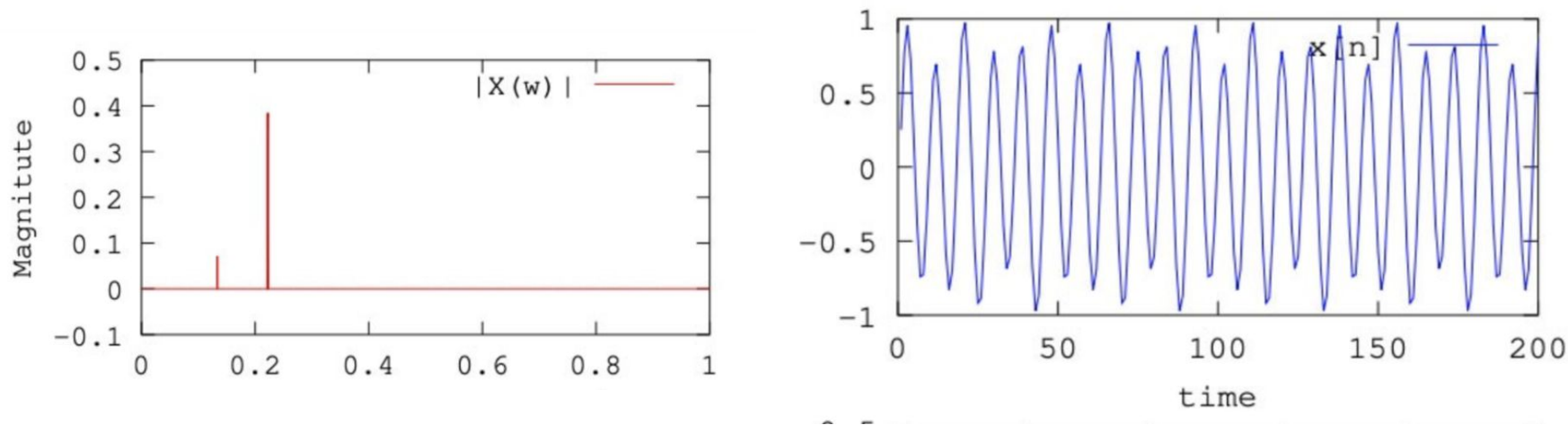


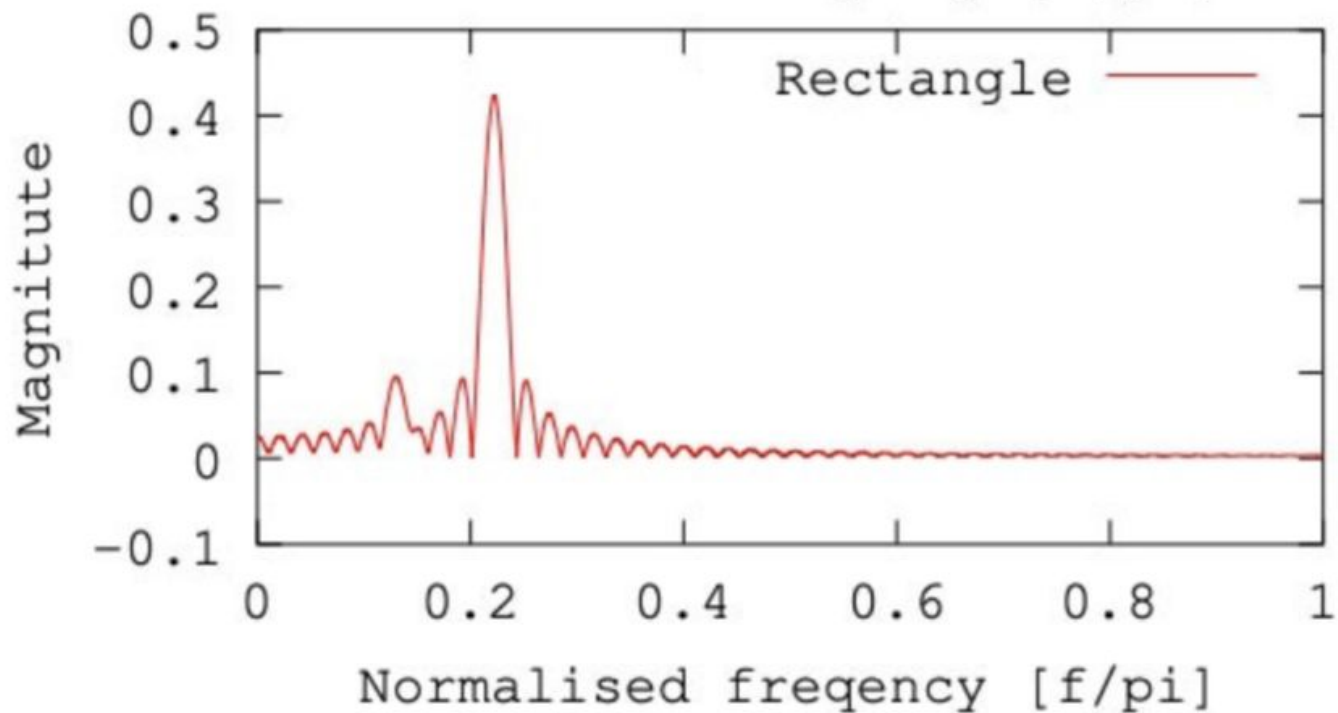
Image from Bryan Pellom

# Windowing example (source)



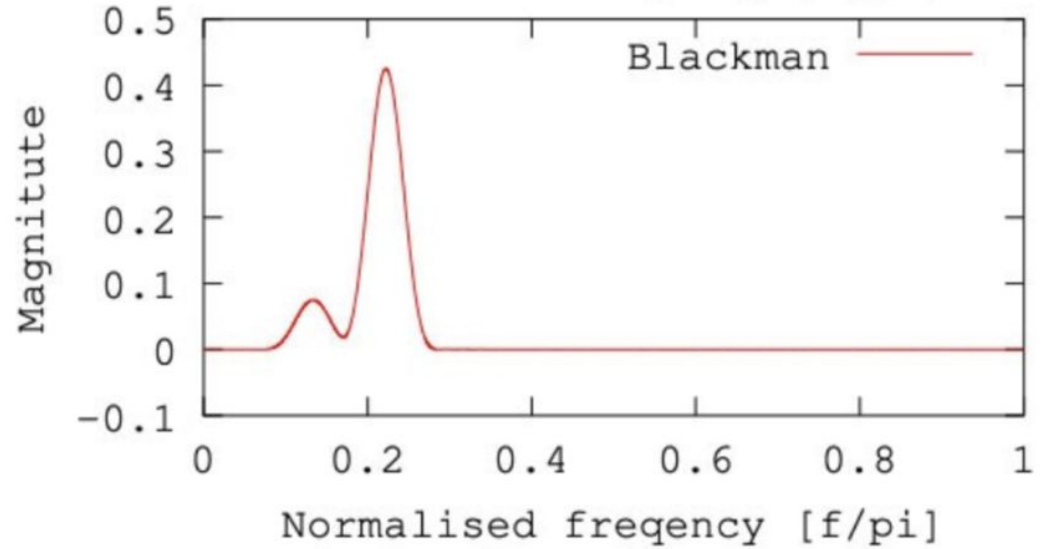
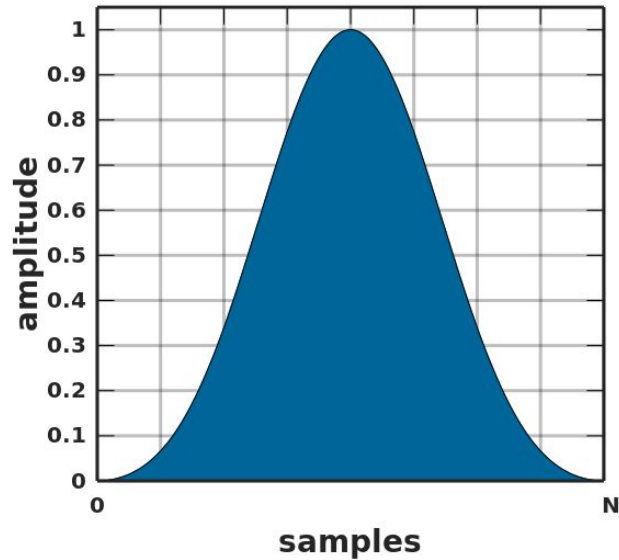
[https://medium.com/@jonathan\\_hui/speech-recognition-feature-extraction-mfcc-plp-5455f5a69dd9](https://medium.com/@jonathan_hui/speech-recognition-feature-extraction-mfcc-plp-5455f5a69dd9)

# Windowing (naive, rectangular result)



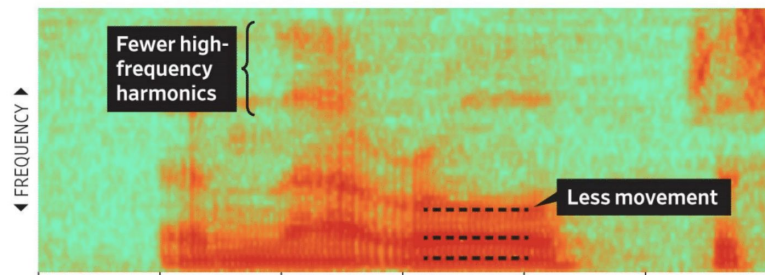
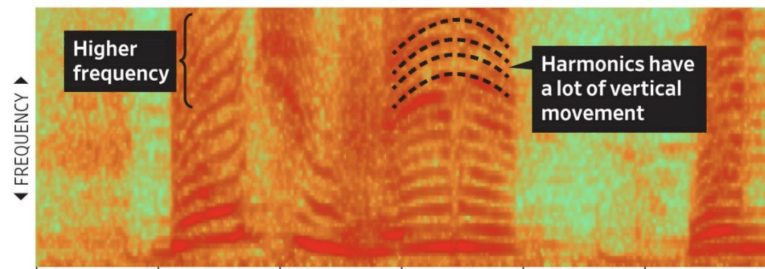
# Windowing (with smoothing)

Blackman window





# And that way WAV turns into a...

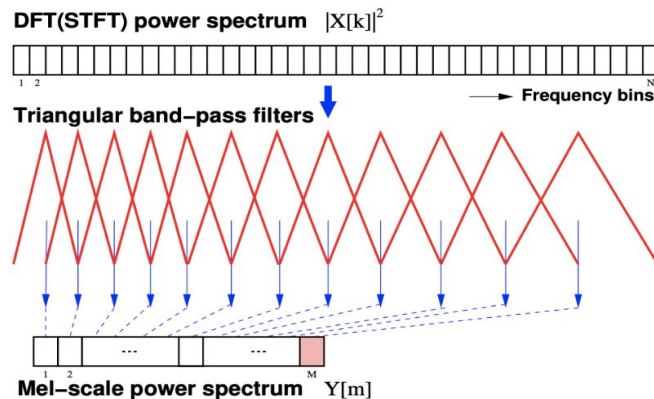


Specific augmentation is usually here too

## Also you can add:

- Pre emphasis ( $x[t]=x[t] - \alpha x[t-1]$ )
- Delta + Delta-Delta
- Normalization

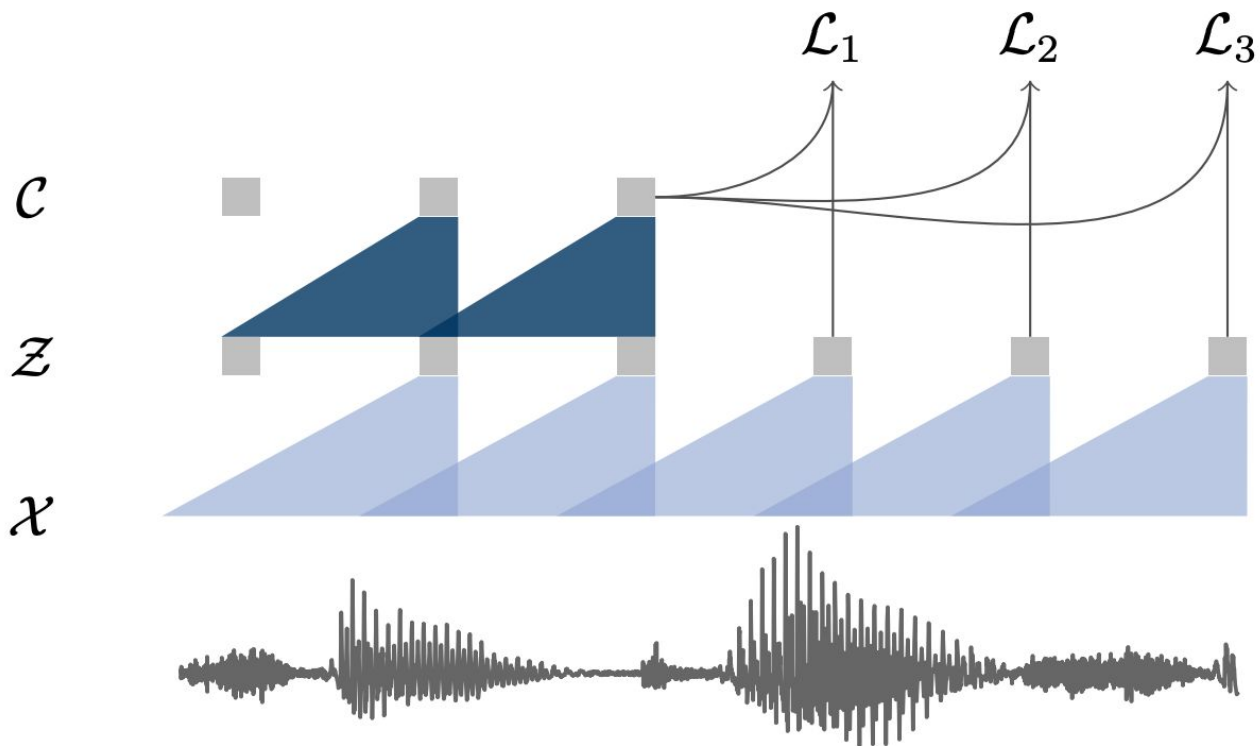
- Mel Scale



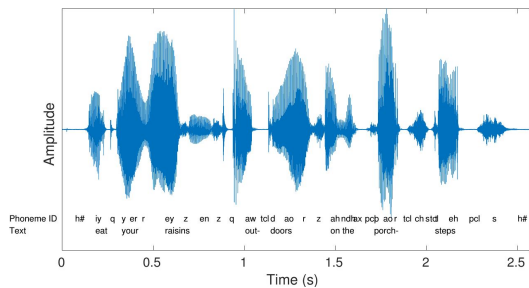
$$Y_t[m] = \sum_{k=1}^N W_m[k] |X_t[k]|^2$$

- Log
- Apply more functions... (IDFT for MFCC etc.)

# Wav2Vec



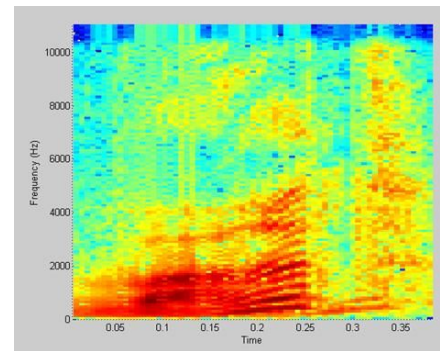
# TL; DR



Sound



**MAGIC BLACK BOX**



Image

# Other “small” tasks to be mentioned:

- VAD (Voice Activity Detector)
- EOU (End of Utterance)
- Spotter
  - Phrase spotter (OK, Google! etc.)
  - Action spotter (a gunshot, a car passing by, etc.)
- Biometry-похожие:
  - Identification (which one of the preset N speakers was that?)
  - Verification (is speaker from the “permitted” set or not?)
  - Speaker features extraction (M-F, child vs adult, age, etc.)
- Acoustic Scene Classification