# Exploring the best way to represent an audio signal for ML model
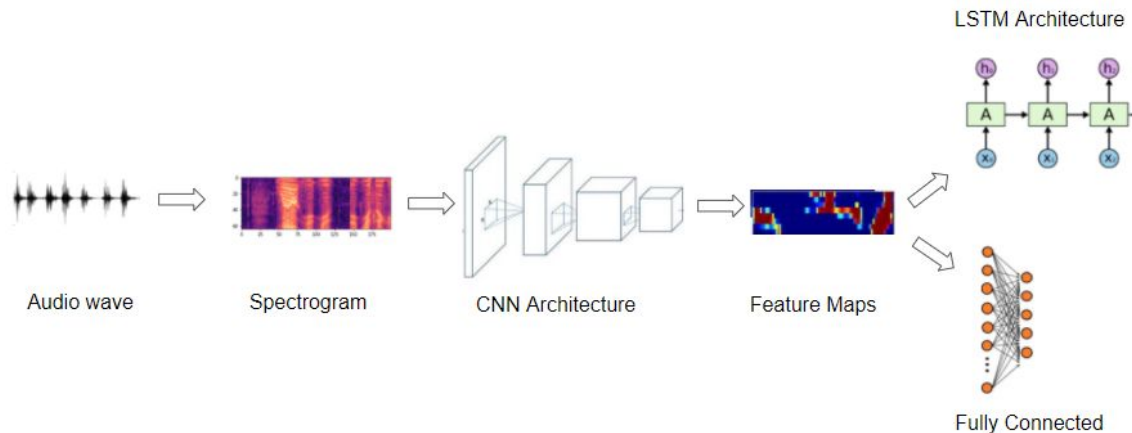
On the example of audio classification
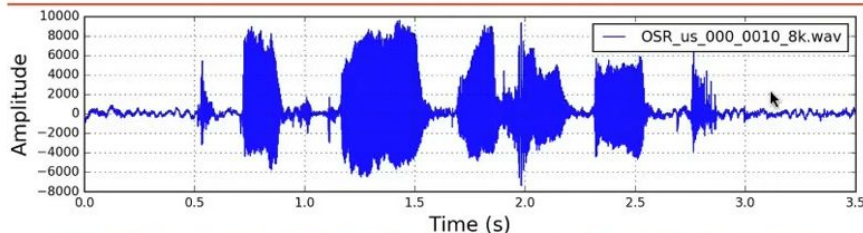
★ Raw signal data
★ Short-Time Fourier-Transform (STFT)
★ Mel-spaced Filterbank
★ Mel Spectrogram
★ Discrete Cosine Transform (DCT)



LSTM Architecture

Audio wave    Spectrogram    CNN Architecture    Feature Maps

Fully Connected

# Looking closely at audio signal post-processing

## Mel-Spectrogram and Mel-Frequency Cepstral Coefficients (MFCCs)

https://haythamfayek.com/2016/04/21/speech-processing-for-machine-learning.html
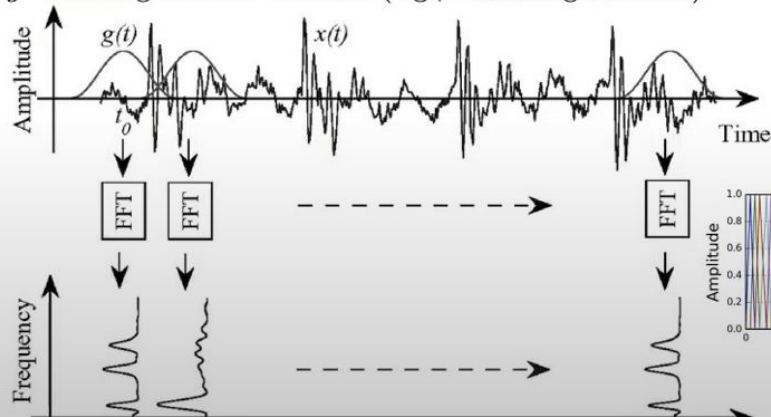
**Short-Time Fourier-Transform (STFT)**

$x \rightarrow$ signal in the time domain (e.g., $x \in \mathbb{R}^{28,200=3.525 \text{ seconds} \times 8,000 \text{ Hz}}$)
  sampling frequency (number of samples per seconds)

$g \rightarrow$ sliding window function (e.g., Hamming function)



$X_i \in \mathbb{R}^{200} \rightarrow i$-th frame of signal $x$ (25ms frames)

$80 \rightarrow$ frame step (10ms) $\implies X \in \mathbb{R}^{200 \times 350}$ ($350 = (28,200 - 200)/80$)

$\tilde{X}_i \in \mathbb{C}^K \rightarrow$ discrete Fourier transform of $X_i \implies \tilde{X} \in \mathbb{C}^{K \times 350}$

$$\tilde{X}_i(k) = \sum_{n=1}^{N} X_i(n)g(n)e^{-j2\pi kn/N}, k = 1, \ldots, K \qquad N = 200$$

$K = 257 \rightarrow$ number of discrete Fourier transform coefficients

$P_i(k) = \dfrac{1}{N}|\tilde{X}_i(k)|^2 \rightarrow$ Periodogram estimate of the power spectrum

$\implies P \in \mathbb{R}^{257 \times 350}$

**Mel-spaced Filterbank**

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

https://www.youtube.com/watch?v=WJI-17MNpdE

$300$ Hz $\rightarrow$ lower frequency

$4,000$ Hz $\rightarrow$ upper frequency

$M(f) = 1,125 \ln(1 + f/700) \rightarrow$ convert frequency to Mel scale

$M^{-1}(m) = 700(\exp(m/1,125) - 1) \rightarrow$ convert Mel scale to Hz

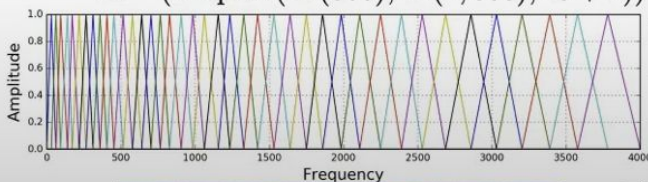$M^{-1}(\text{linspace}(M(300), M(4,000), 26 + 2))$ $\qquad 26 \rightarrow$ number of triangular filters

$T \in \mathbb{R}^{26 \times 257} \implies E = TP \in \mathbb{R}^{26 \times 350}$

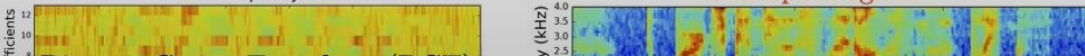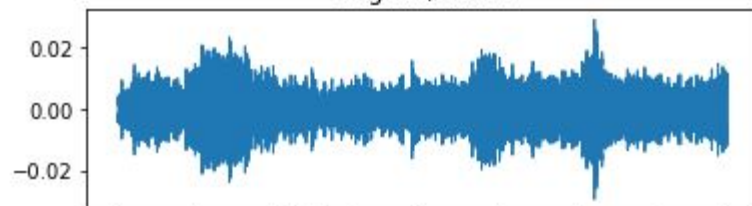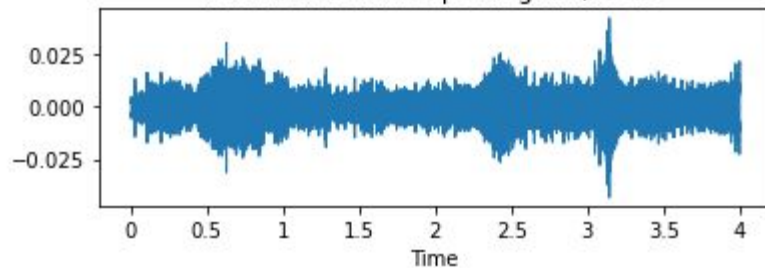$E_i(l) \rightarrow$ amount of energy in filter bank $l$ at frame $i$

$\log(E) \in \mathbb{R}^{26 \times 350} \rightarrow$ log filter bank energy
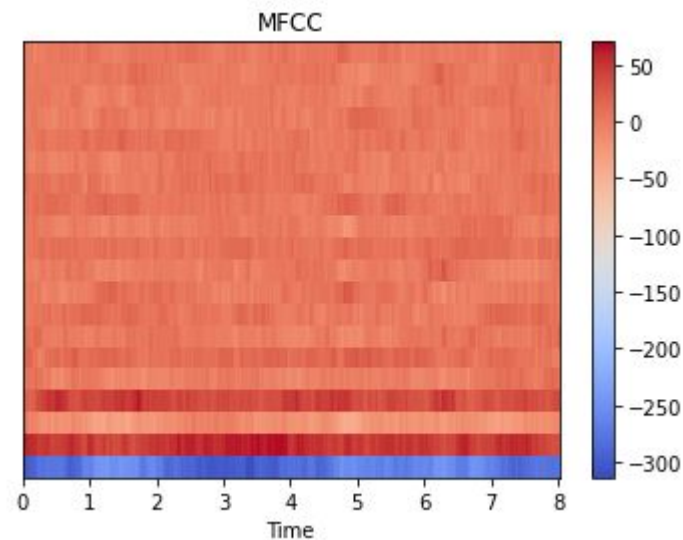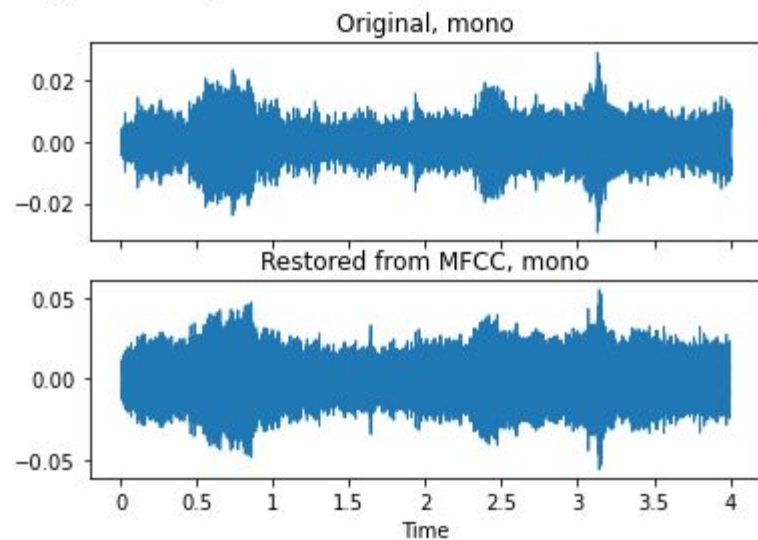
**Mel Spectrogram**

# Mel-Spectrogram analysis

# MFCC analysis

original memory use = 176400, MFCC = 6900

# Mini-Conclusion, based on applying methods to signal

- If one wants to use MFCC as a input for neural network against Mel-Spectrogram, s/he wins approximately **4x** in space (usually it is an input for the first convolutional layer)
- Nobody uses raw data for the input in the latest researches, as it seems hard for a neural network extract features from such an input
- each of methods have information lost after its application

# Evaluating feature extraction of two methods: MFCC vs Mel-Spectrogram – Results

For a paticular Urban8k classification task MFCC performs better (0.77 > 0.65) (accuracy scores)

It is useful to use MFCC instead of Mel-Spectrogram for a Neural Network to be able to increase the batch size if one wants a fast training

# Research conclusion

- Mel-Spectrogram is commonly used for NLP, Transformers
- MFCC is commonly used for Reinforcement Learning, Hidden Markov Processes (when we care about memory)
- It is better to try both methods of pre-processing data for Neural Network, and then take the better.

Possible ways of future research:

    a. As images can be represented, as signals, try to use Mel-Spectrogram or MFCC for its feature extraction

    b. Train more specific architectures to explore more applications of these methods