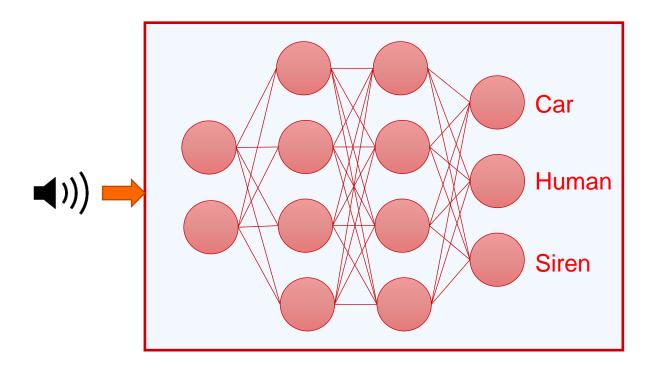


Audio Analytics





Reference:

- Müller, Meinard. Fundamentals of Music Processing: Using Python and Jupyter Notebooks. Vol. 2. Cham: Springer, 2021.
- Knees, Peter, and Markus Schedl. Music similarity and retrieval: an introduction to audio-and web-based strategies. Vol. 36. Heidelberg: Springer, 2016.
- Bäckström, Tom, Okko Räsänen, Abraham Zewoudie, Pablo Perez Zarazaga, and Sneha Das. "Introduction to speech processing." (2020). https://speechprocessingbook.aalto.fi/
- https://librosa.org/doc/main/index.html
- https://github.com/musikalkemist/AudioSignalProcessingForML/tree/master

Why Audio Analytics on Edge





- Low latency is critical for fast action (e.g., audio scene classification)
- Sensitivity/Security of the data is critical for many applications
- Redundancy of the data can be reduced (enable compression) for further analysis (offline)

Audio Analytics on Edge - Applications

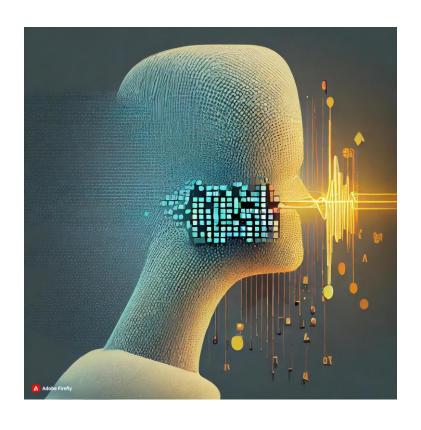




- Audio Scene Classification
 - e.g. surveillance
- Audio Event Detection
 - e.g. screeching brakes
- Speaker/Speech Recognition
 - Keyword Recognition
- Speech to Text (or vice versa)
- Conversational User Interfaces (CUI)

Content





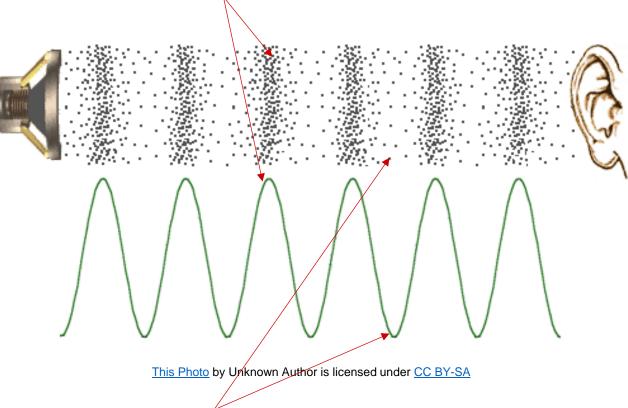
- Sound Waves and Digitization
- Time-domain and frequency-domain features
- Audio transformations
- Lightweight Machine Learning for Audio

Introduction to Sound Waves





Rarefaction

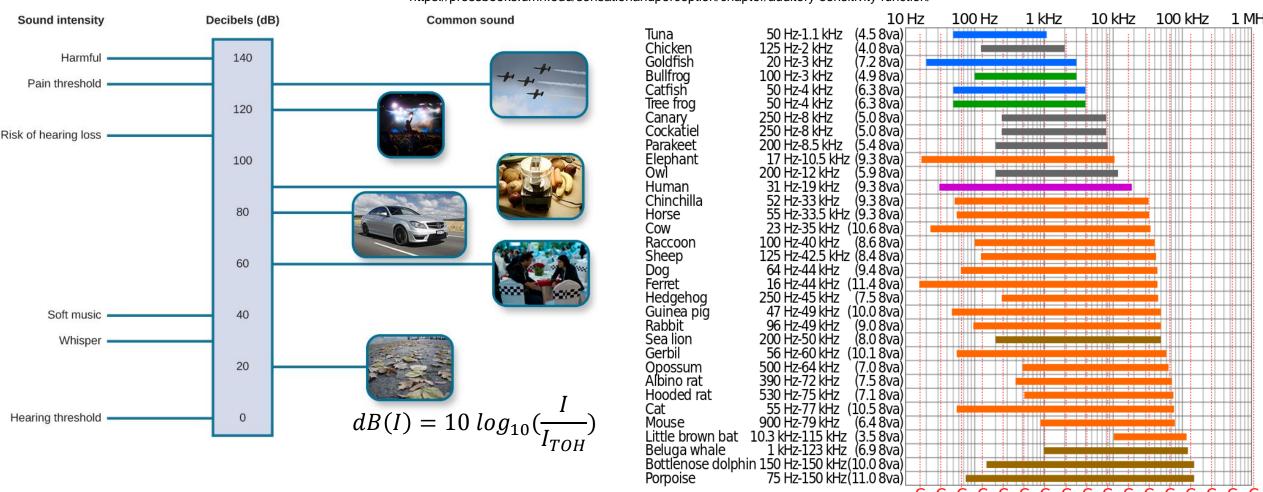


 $y(t) = A \sin(2\pi f t + \varphi)$ A \rightarrow Amplitude
f \rightarrow frequency
t \rightarrow time $\varphi \rightarrow$ phase

Sound Amplitude and Frequency



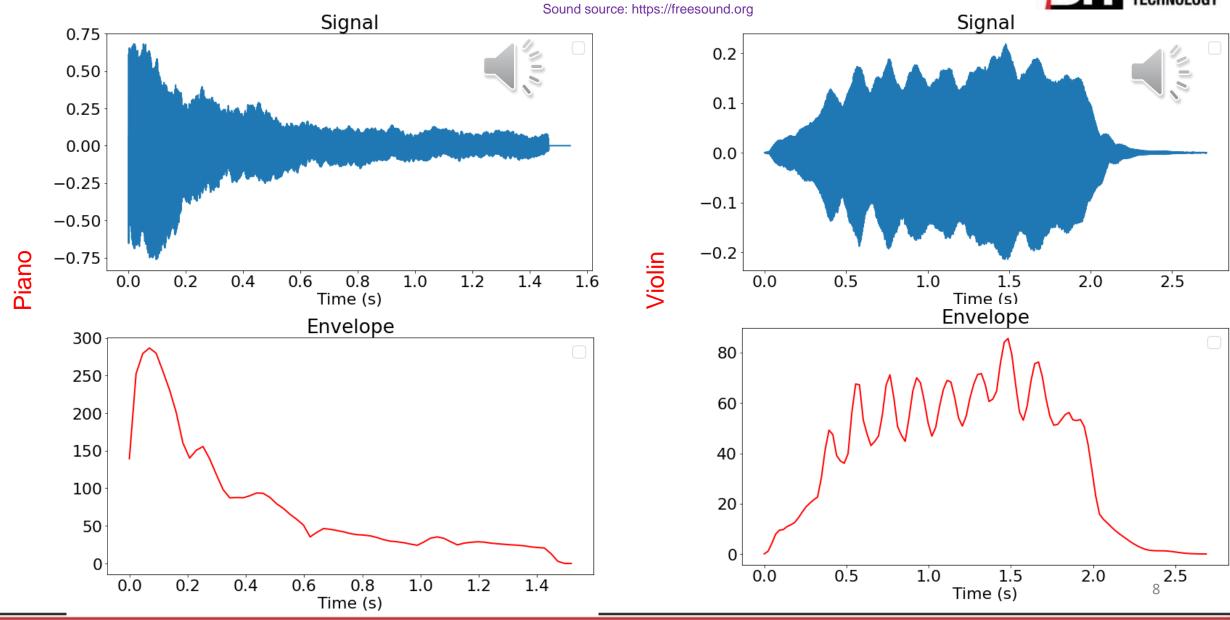
https://pressbooks.umn.edu/sensationandperception/chapter/auditory-sensitivity-function/



Simple machine learning is to put thresholds on amplitude and frequency (but world is noisy!)

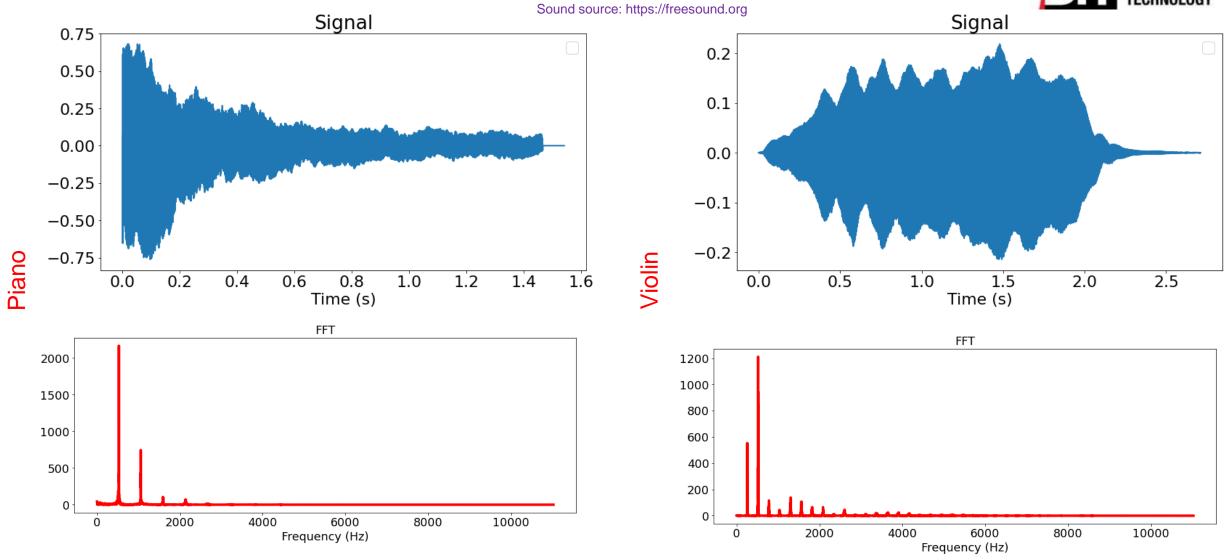
Sound Amplitude (Envelope)





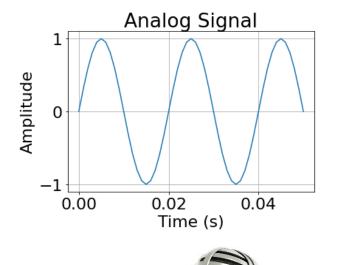
Sound Amplitude (Frequency)

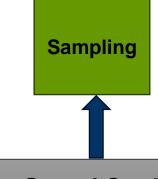




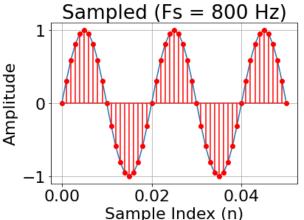
Digitization of Sound Waves

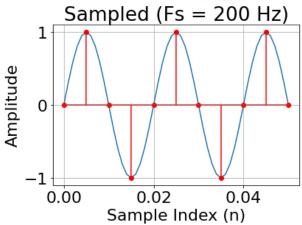




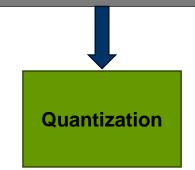




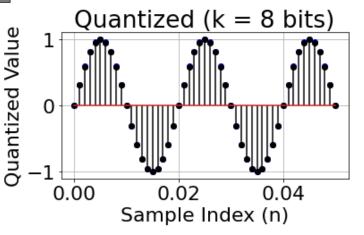


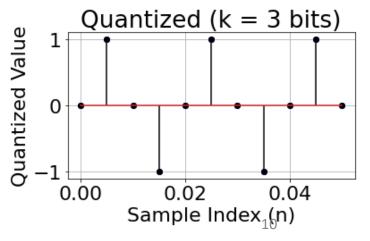






- Sampling frequency > 2 times the max frequency (bandwidth)
 - Nyquist Criteria
 - Sound cards typically support up to 48 kHz.
- Higher bits (resolution) is better!

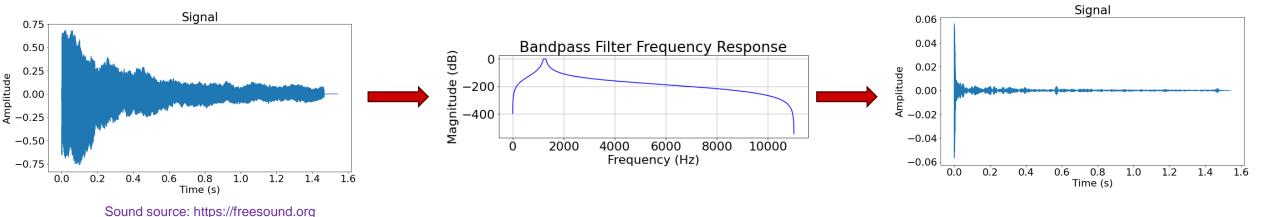




Pre-processing



- Noise reduction: background, artefacts
- Standardization: formats, sampling rate, varying lengths
- Normalization: varying signal amplitudes
- Improved learning models: efficient training, low complexity



Features from Sound





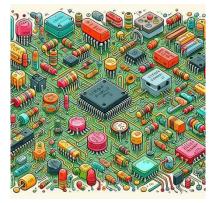
High-Level

Examples: Melody, rhythm, tempo, genre etc.



Mid-Level

Examples: Pitch and beat related, Spectrogram, Mel-Spectrograms, MFCCs



Low-Level

Examples: Amplitude envelope, energy, spectral centroid, zero crossing rate

Time Domain Features

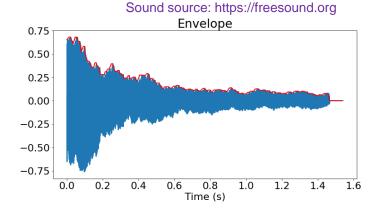


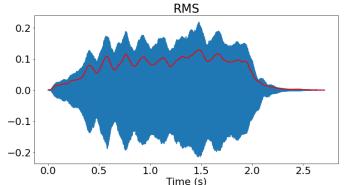
Envelope, $AE_k = max(s(k)) for k - K/2: k + K/2$

- Divide the total time into small time windows
- Find the max/min/both in the window
- Useful in onset detection, music genre classification
- Very sensitive to outliers!
- Sensitive to noise!

Root Mean Square,
$$RMS_k = \sqrt{\frac{1}{K} \sum_{k=K/2}^{k+K/2} s(k)^2}$$

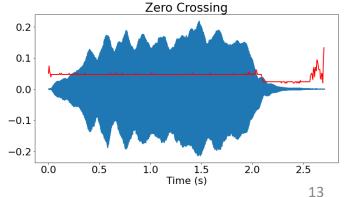
- Divide the total time into small time windows
- Find the root mean square in the window
- Less sensitive to outliers
- Audio segmentation, music genre classification
- Sensitive to noise!





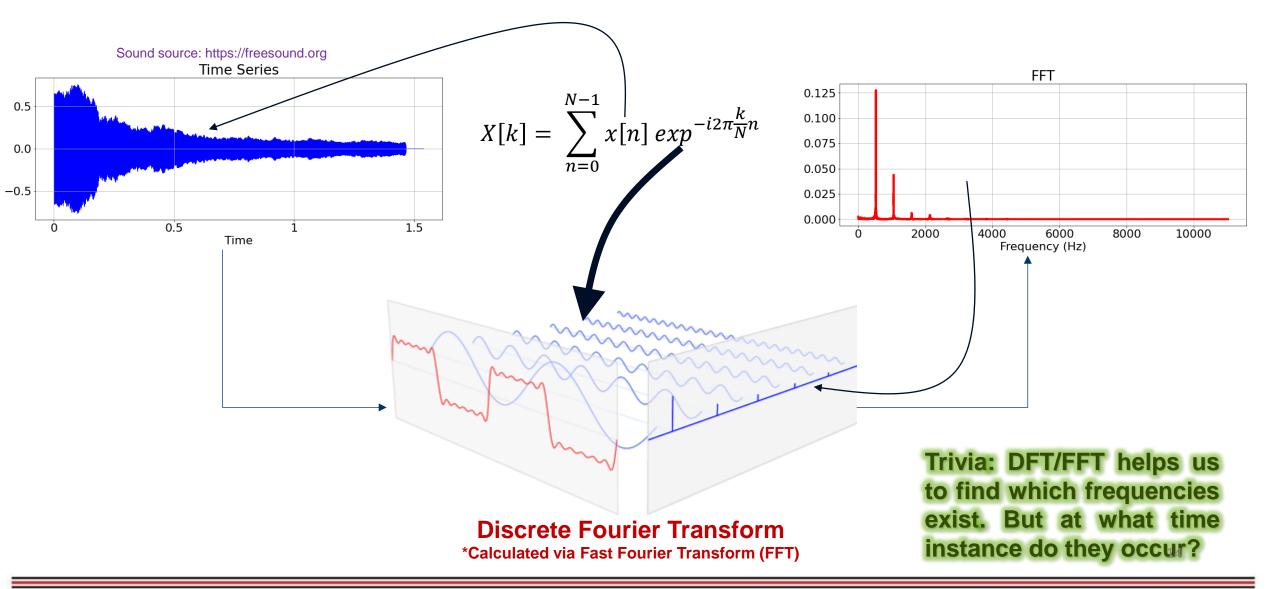
Zero Crossing Rate, $ZCR_k = \frac{1}{2} \sum_{k=K/2}^{k+K/2} |sgn(s(k)) - sgn(s(k+1))|$

- Number of times a signal crosses the horizontal axis
- Music genre recognition, pitch estimation, speech segmentation
- Sensitive to noise!



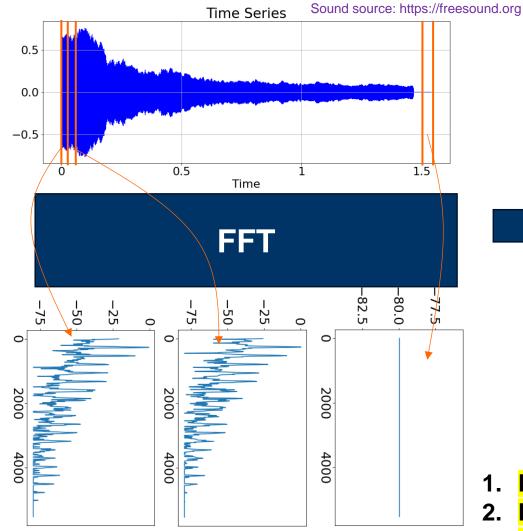
Frequency Domain

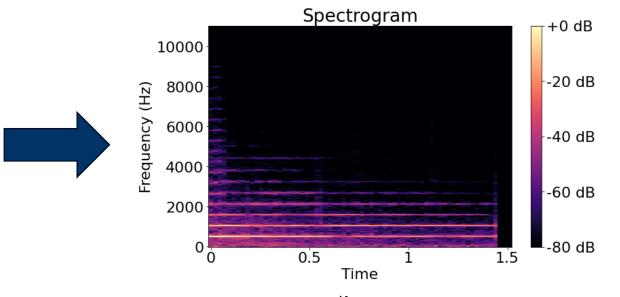




Short Time Fourier Transform - Spectrograms







$$S(k,m) = 20 \log 10 \left(\left| \sum_{n=0}^{N_W - 1} x[n] w[m-n] exp^{-i2\pi \frac{k}{N}n} \right| \right)$$

- 1. Divide the total time into small time windows
- 2. Find FFT for each window (after multiplying with a short decaying windows (e.g. Gaussian)
- Stack all the FFTs (actually FFT power) to plot an image
- 4. Take log of FFTs to show small changes (result will be dB)

Mel - Spectrograms



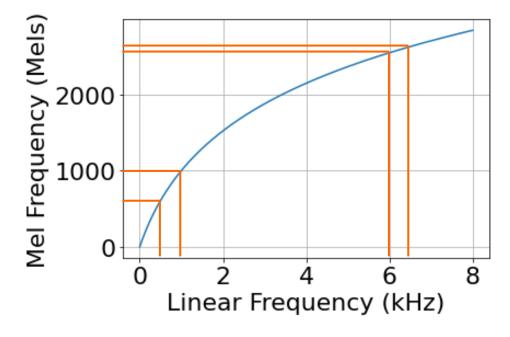
Humans perceive sound in log frequency scale and not linear frequency scale

The mel-scale is "a perceptual scale of pitches judged by listeners to be equal in distance from one another"





$$mel_f = 2595 \ log_{10} \left(1 + \frac{lin_f}{700} \right)$$
(O'Shaughnessy 1987)



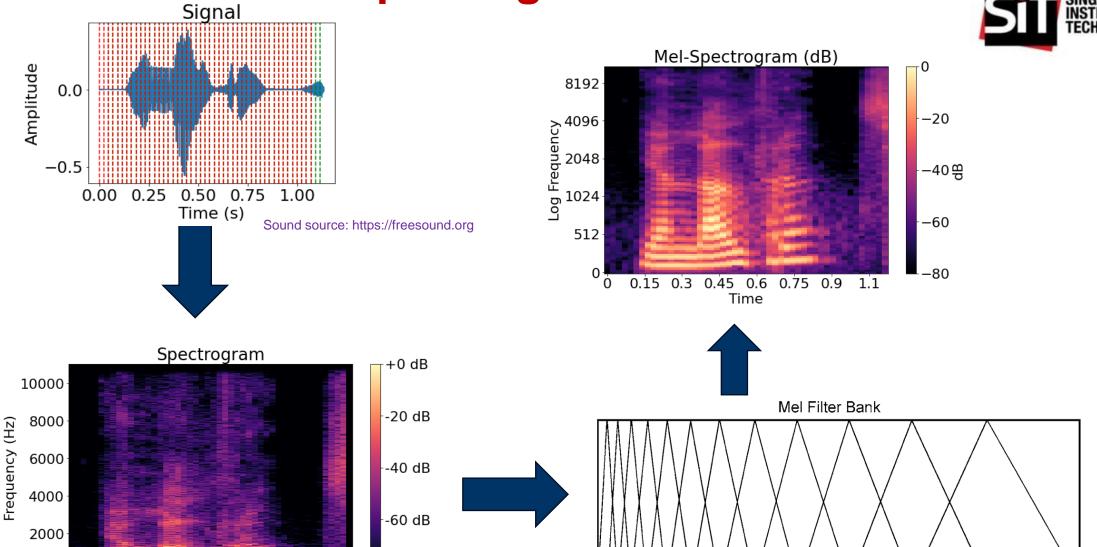
Mel - Spectrograms

-80 dB

0.15 0.3 0.45 0.6 0.75 0.9

Time

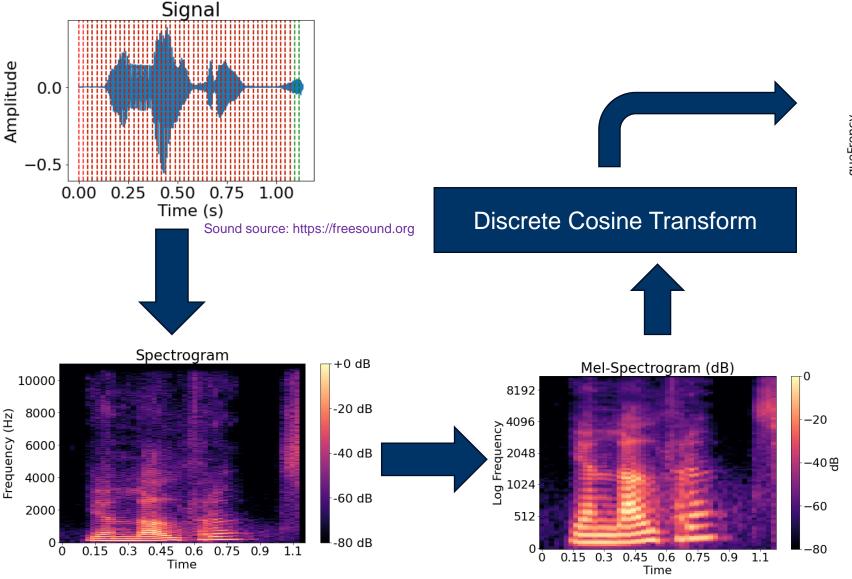


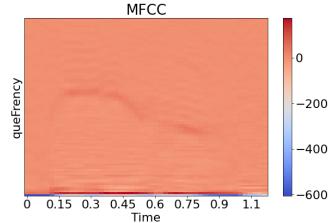


Frequency (Hz)

Mel Frequency Cepstral Coefficients (MFCC)







- A spectrum of spectrum is termed as "Cepstrum"
- MFCC is in time-domain and is much less decluttered and hence useful for Machine Learning tasks
- MFCC is considered as a representation of timbre (perceived sound quality)

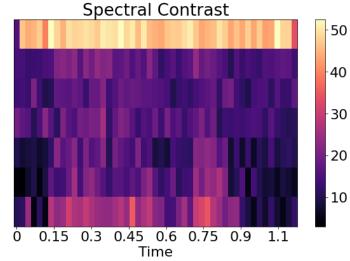
18

Frequency Domain Features



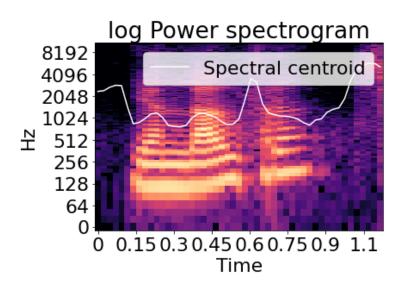
Spectral Contrast,
$$BER[n] = \frac{\sum_{k=0}^{F-1} S_n(n)^2}{\sum_{k=F}^{N-1} S_n(n)^2}$$

- Comparison of energy in the lower/higher frequency bands
- Music/speech discrimination



Spectral Centroid, $SC_n = \frac{\sum_{k=0}^{N-1} S_n(n) n}{\sum_{k=0}^{N-1} S_n(n)}$

- Centre of gravity of magnitude spectrum
- Frequency band where most of the energy is concentrated
- Measure of brightness of sound!
- Classification



Low Cost (Compute) Audio Analytics



Amplitude Envelope **Zero Crossing Rate** RMSE

Time Domain

Spectral Centroid Spectral Contrast Frequency Domain

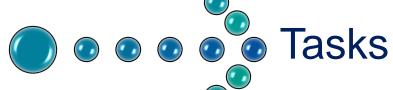
Machine Learning











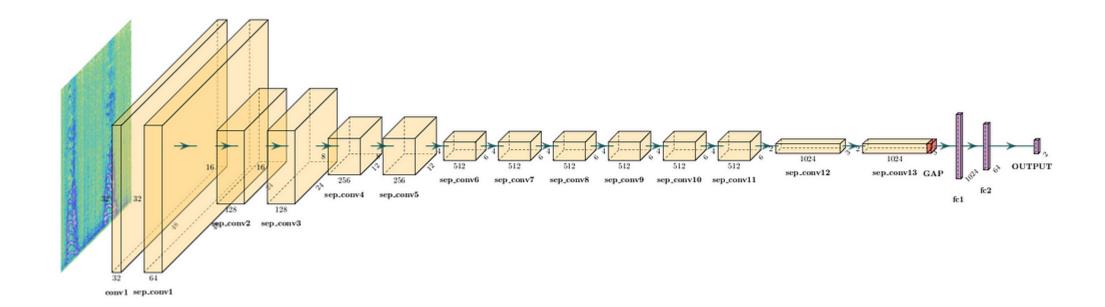
Spectrogram MFCC

Time-Frequency

State-of-the art



YAMNet: A deep convolutional neural network that predicts 521 audio event classes from the AudioSet-YouTube corpus it was trained on. The architecture is based on the light weight MobileNet v1.

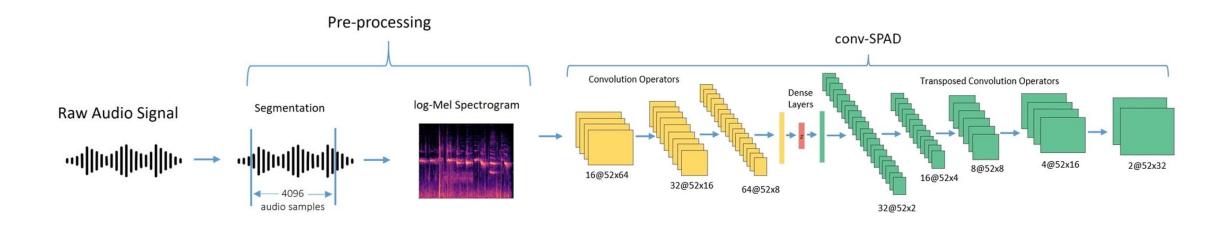


https://www.tensorflow.org/hub/tutorials/yamnet

State-of-the art



conv-SPAD: A simple custom modelled convolutional SPectral audio-based anomaly detection.



Lo Scudo, Fabrizio, Ettore Ritacco, Luciano Caroprese, and Giuseppe Manco. "Audio-based anomaly detection on edge devices via self-supervision and spectral analysis." Journal of Intelligent Information Systems (2023): 1-29.

Build Your Own Model

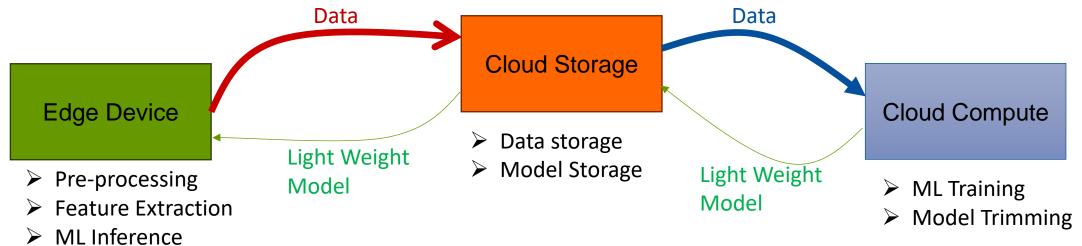


- Pre-processing
 - Normalization, Segmentation
- Feature extraction
 - Spectrogram, MFCC
- Model selection
 - CNNs, RNNs, Transformers
- Transfer learning
 - Leverage existing models
- Useful datasets
 - TAU Urban Acoustic Scenes, TUT Rare sound events, LibriSpeech and Audiset

Lo Scudo, Fabrizio, Ettore Ritacco, Luciano Caroprese, and Giuseppe Manco. "Audio-based anomaly detection on edge devices via self-supervision and spectral analysis." Journal of Intelligental Information Systems (2023): 1-29.

Current Trend – Cloud to Edge for Scalability





- Edge Impulse
- Speech Processing with AI
 - Murf.ai,
 - AWS Transcribe and Polly,
 - Google Cloud Text-to-Speech,
 - Azure Text to Speech API,
 - IBM Watson Text to Speech

Summary





- Audio analytics on edge helps in faster response, data security and enhanced analysis
- Pre-processing of audio is important
- Selection of right input features and model is important
- Build in the cloud and adapt to edge for scalability



Demo of audio analytics @ https://edgeaudioanalytics.streamlit.app/

Contact: mahesh.panicker@singaporetech.edu.sg