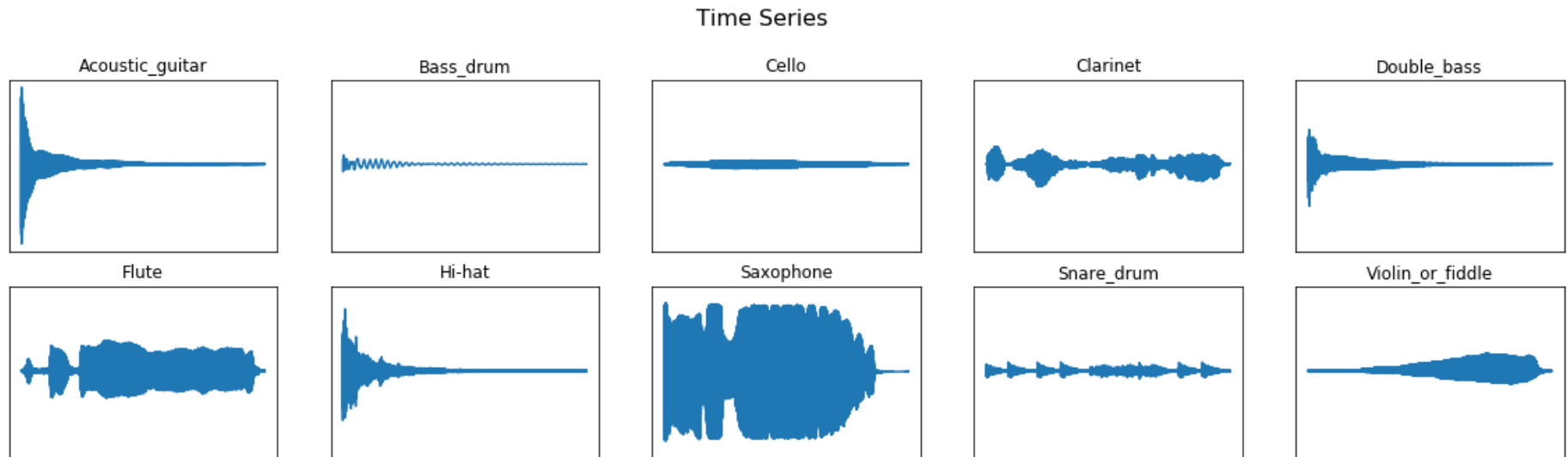

Audio Classification: Learning the Mel Scale

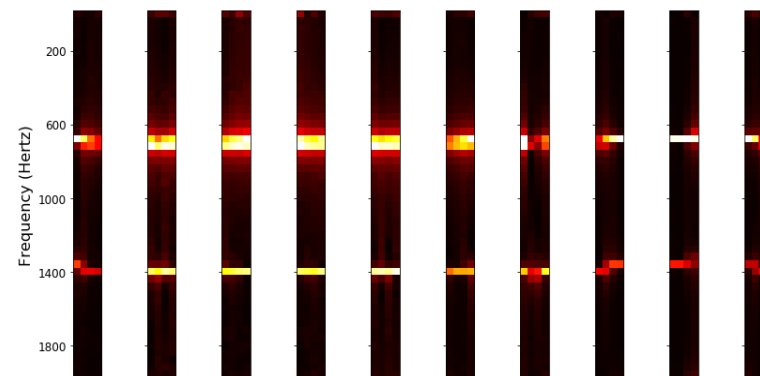
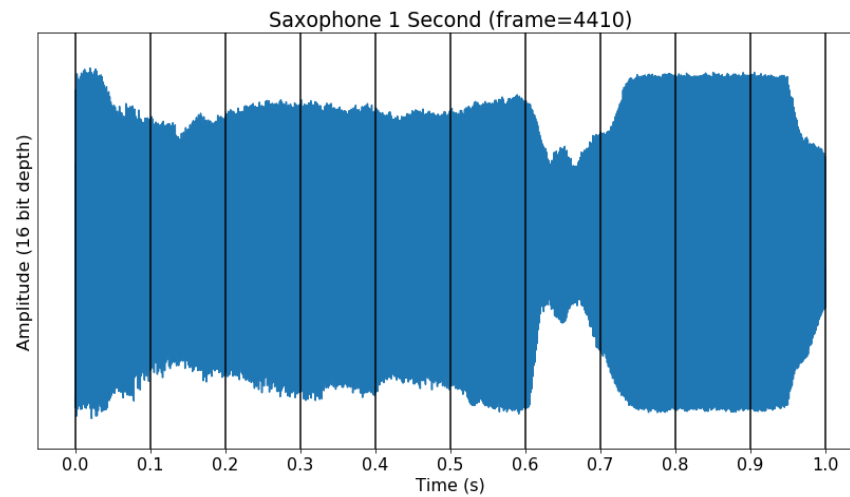
Amitangshu Pal

Audio Classification

- ❑ There are a bunch of sounds coming from different musical instruments
- ❑ Can your smart devices (say alexa or your smartphone) distinguish between them?
 - ❑ Human can do this distinction pretty accurately
 - ❑ How can we include this feature in smart devices?



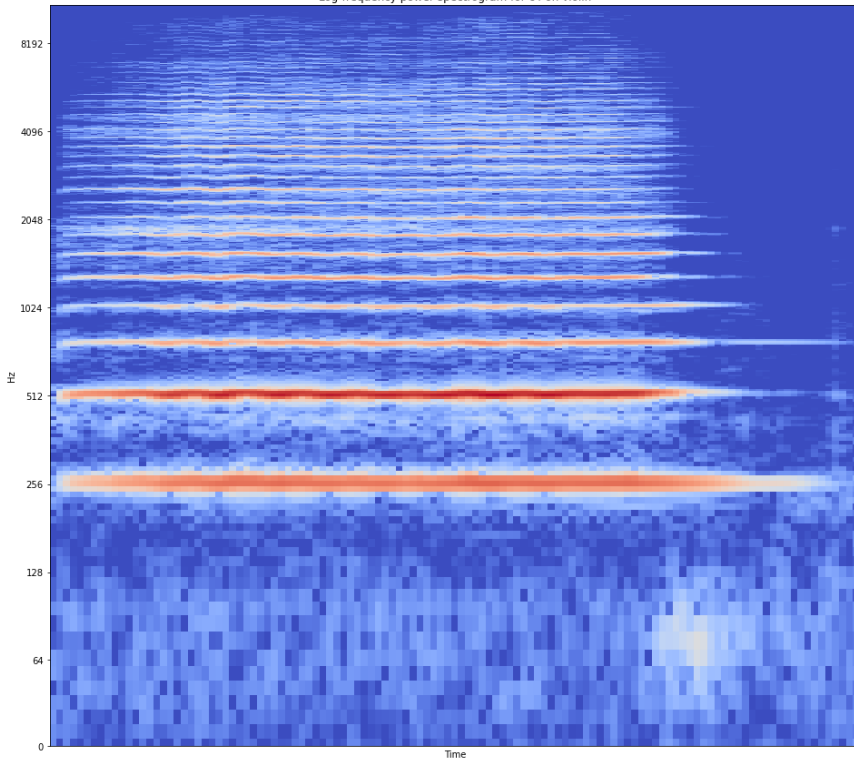
Spectrogram



Spectrogram

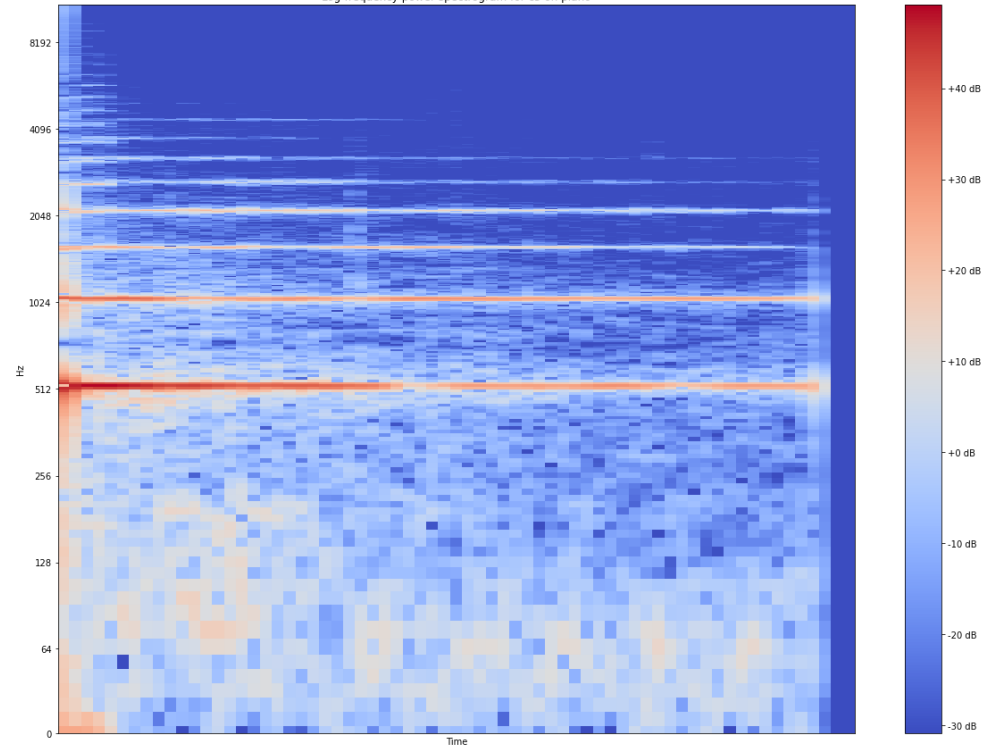


Log-frequency power spectrogram for c4 on violin



Violin

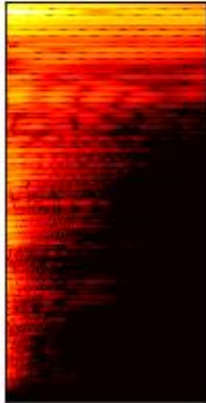
Log-frequency power spectrogram for c5 on piano



Piano

Spectrogram

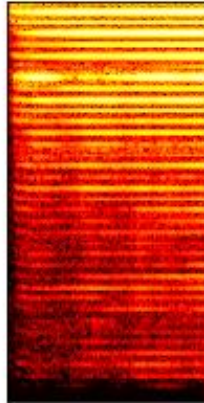
Acoustic_guitar



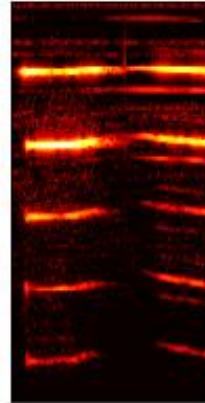
Bass_drum



Cello



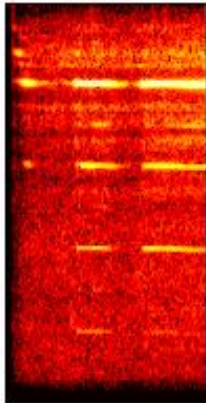
Clarinet



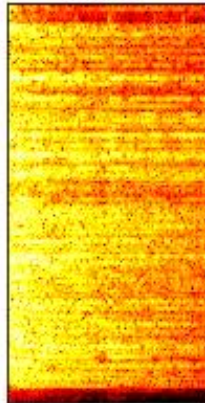
Double_bass



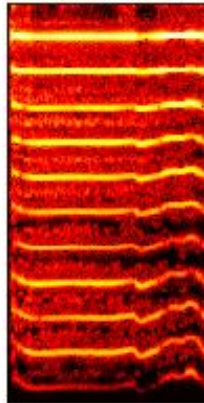
Flute



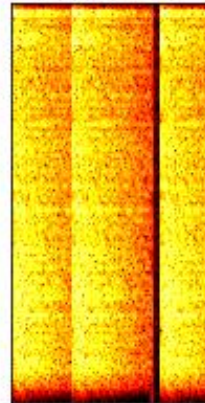
Hi-hat



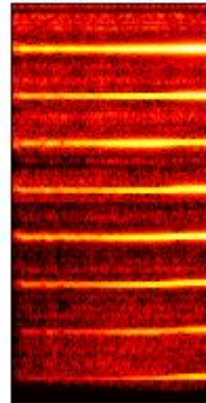
Saxophone



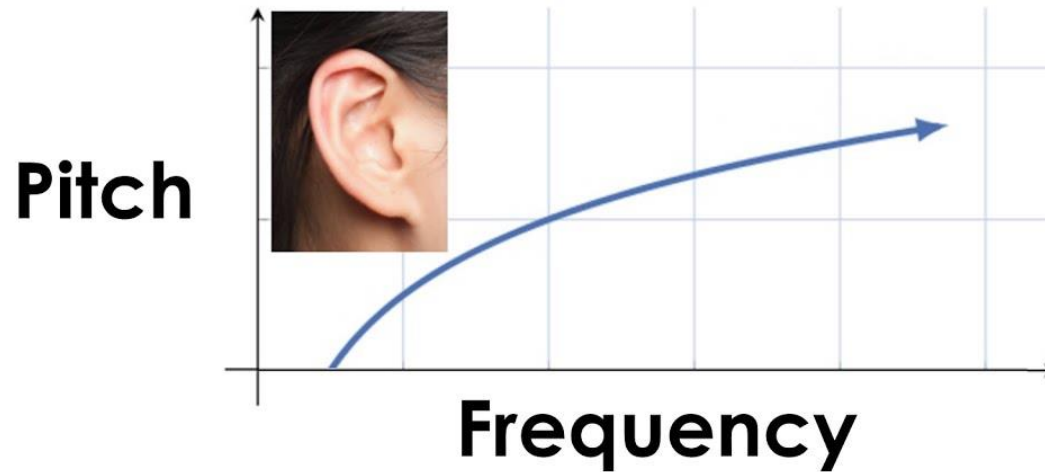
Snare_drum



Violin_or_fiddle



Logarithmic Perception of Frequency



Perception of frequency

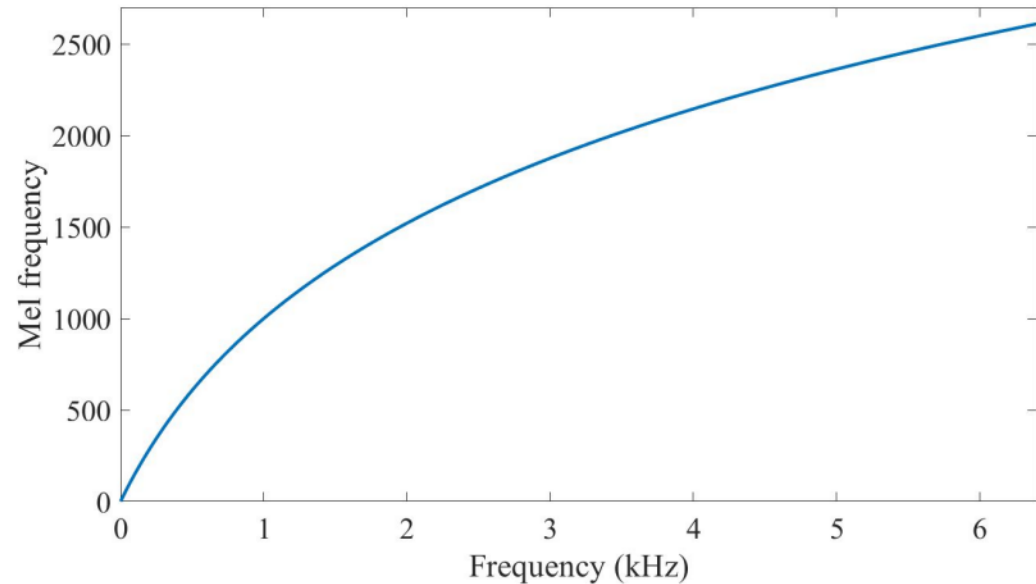
Human perceives frequencies logarithmically

Mel Scale

- Logarithmic scale → equal distances have same **perceptual** distance
- 1000 Hz = 1000 Mel

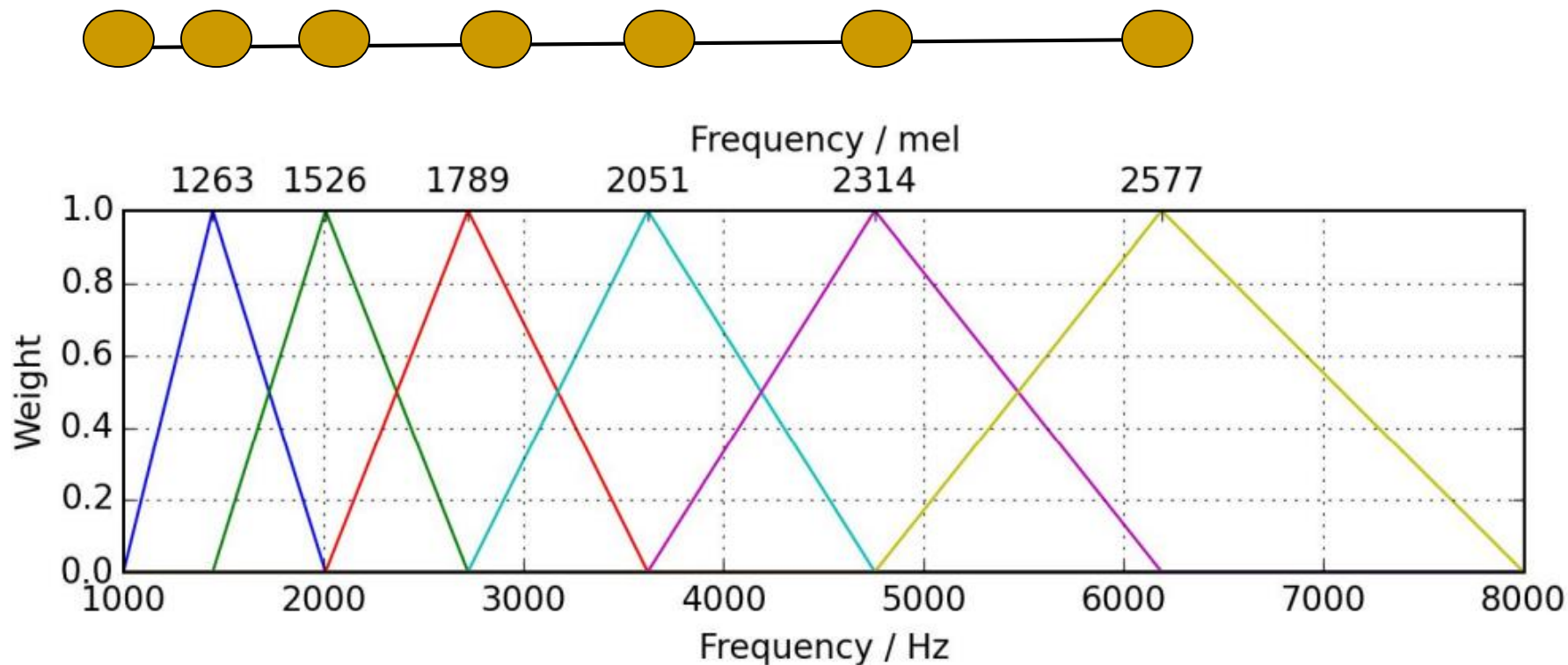
$$m = 2595 \log \left(1 + \frac{f}{700} \right)$$

$$f = 700 \left(10^{\frac{m}{2595}} - 1 \right)$$

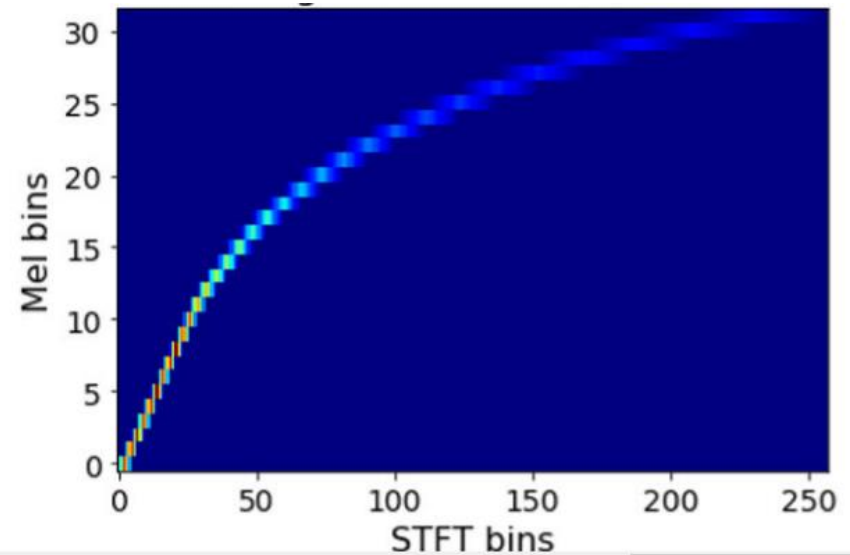
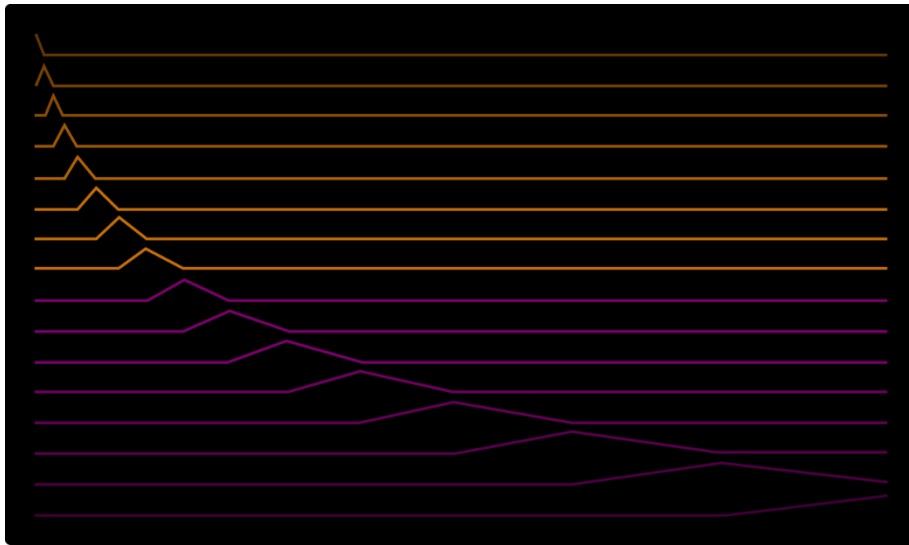


Mel Filter Banks

$$f = 700 \left(10^{\frac{m}{2595}} - 1 \right)$$



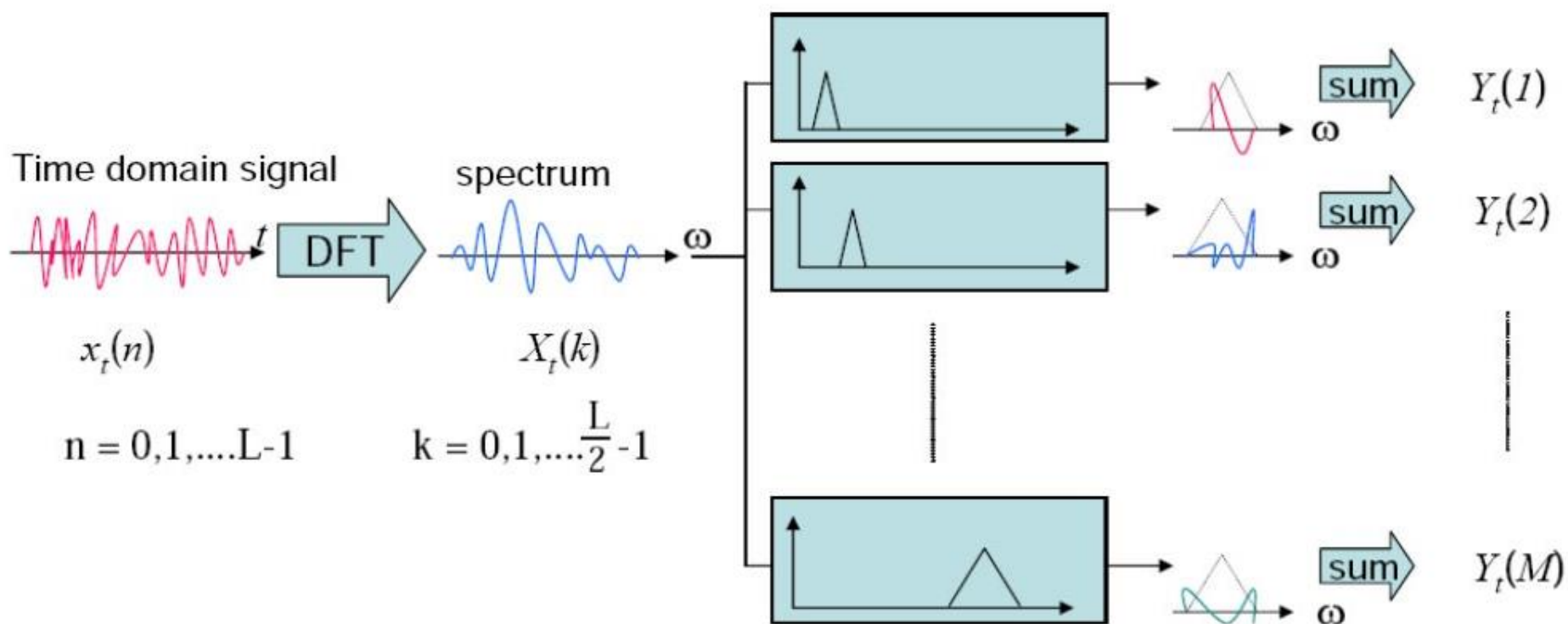
Mel Filter Banks



https://developer.apple.com/documentation/accelerate/computing_the_mel_spectrum_using_linear_algebra

<https://ieeexplore.ieee.org/document/9174990?denied=>

Mel Spectrogram Steps



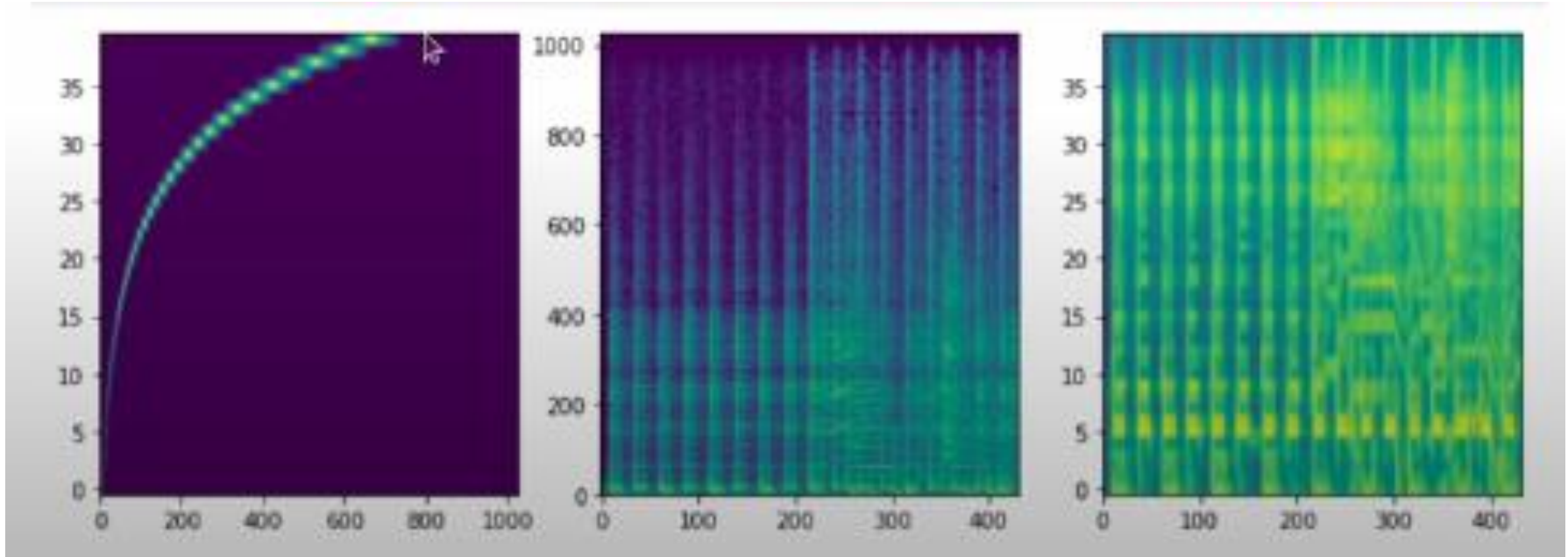
Mel Filter Banks

Filter bank
matrix

Spectrogram
matrix

=

Mel Spectrogram
matrix



40×1024

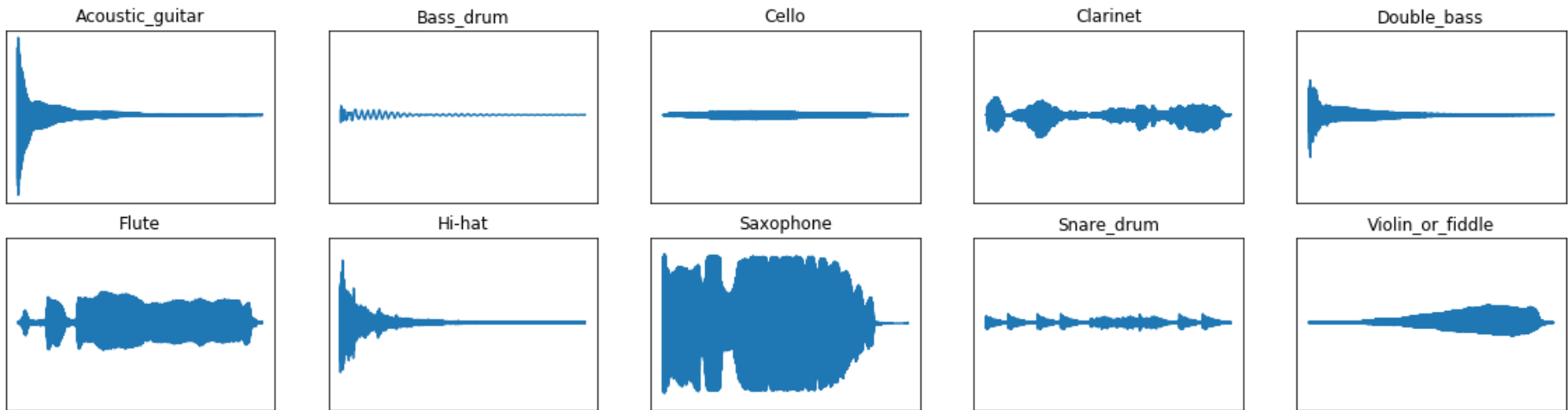
1024×450

40×450

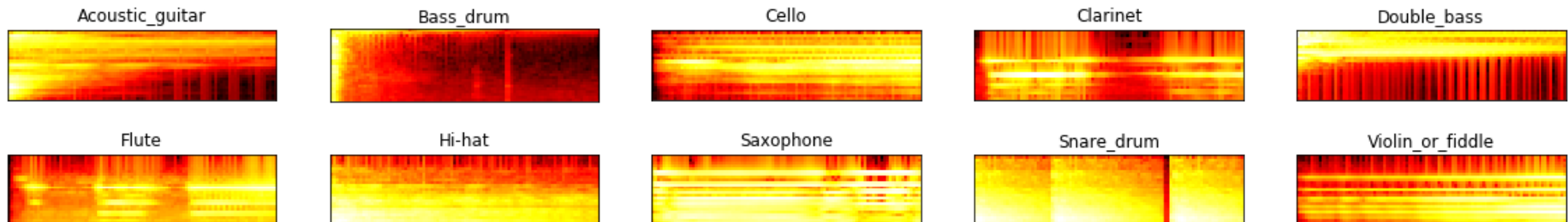
Also can be looked as **dimension reduction** or **lossy compression** of spectrogram

Mel Spectrogram

Time Series



Filter Bank Coefficients



Cepstrum

MFCC → Mel Frequency **Cepstral** Coefficients

Time-domain signal

Spectrum

Log spectrum

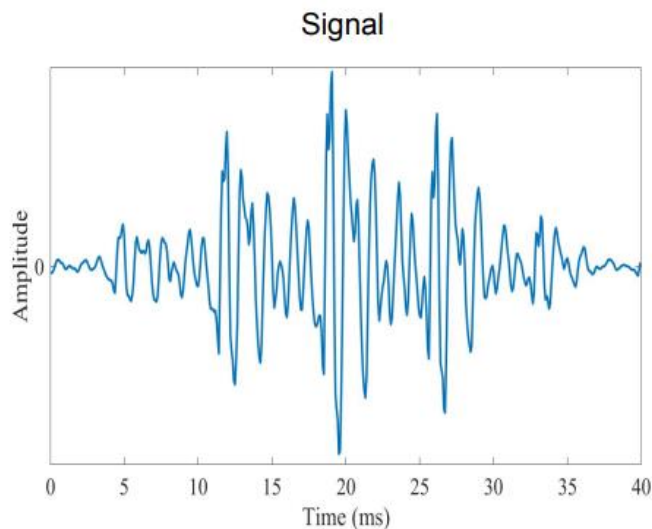
Cepstrum

$$C(x(t)) = F^{-1}[\log(F[x(t)])]$$

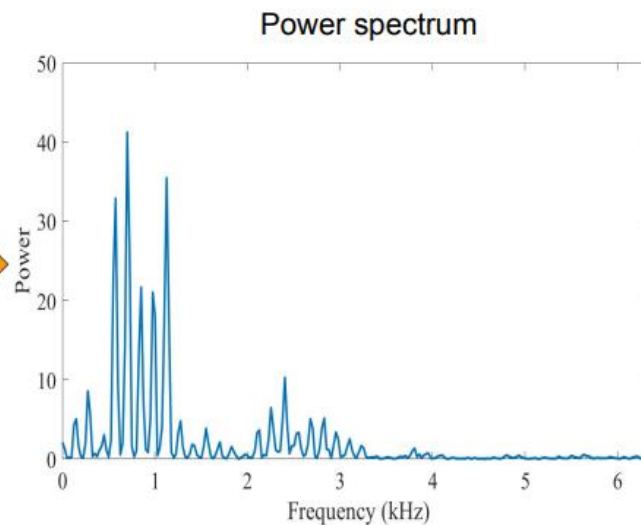
Cepstrum

$$C(x(t)) = F^{-1}[\log(F[x(t)])]$$

Time-domain signal (pink box)
Spectrum (blue box)
Log spectrum (orange box)
Cepstrum (green box)

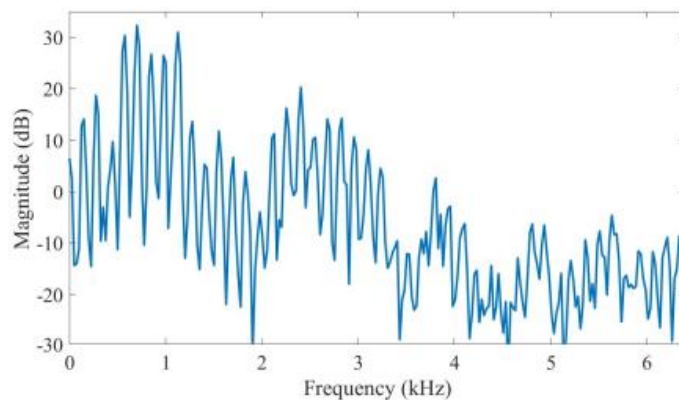


DFT



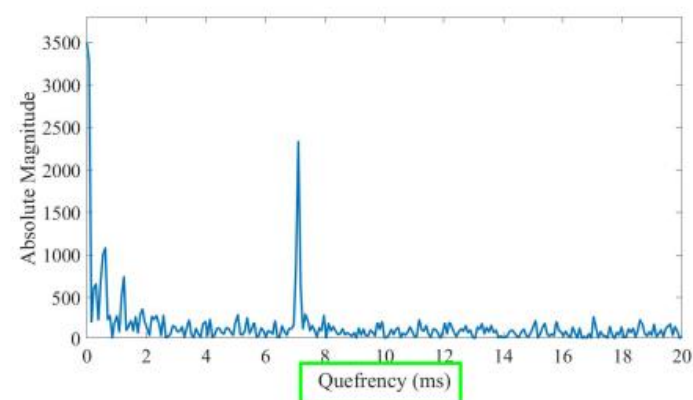
Log power spectrum

Log



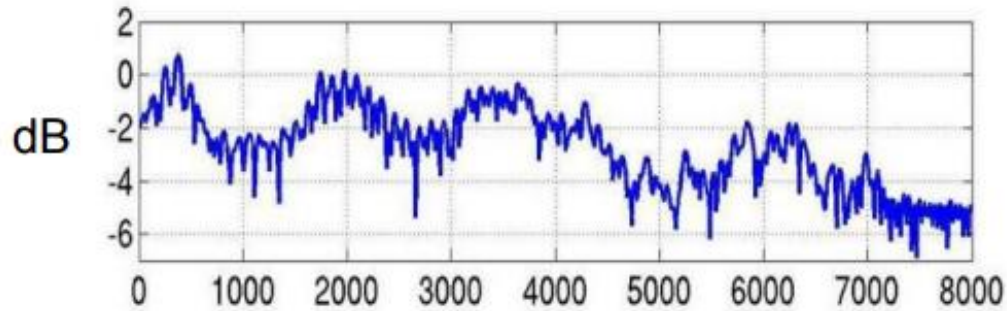
IDFT

Cepstrum

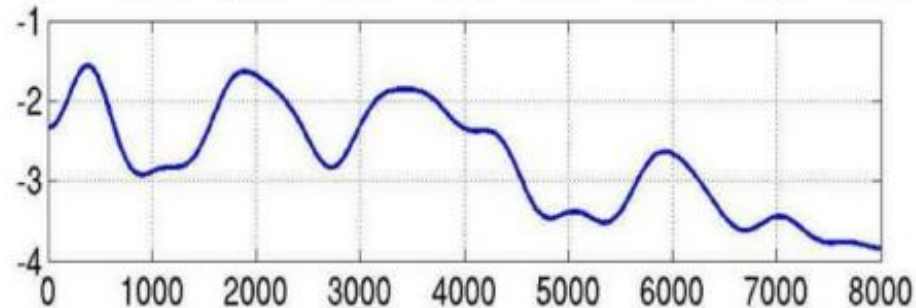


Understanding Cepstrum

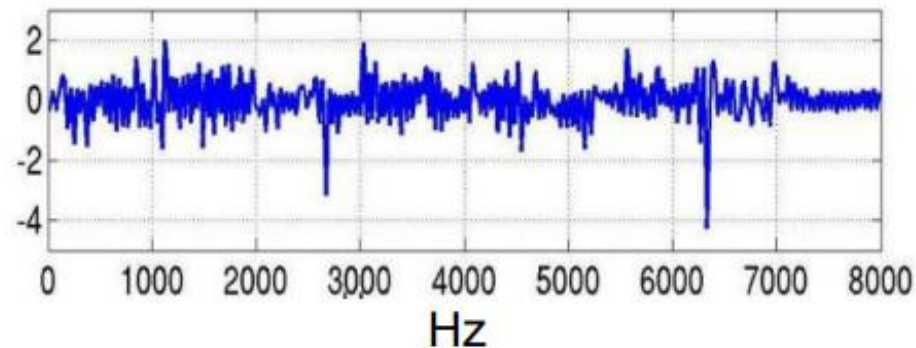
Log-spectrum



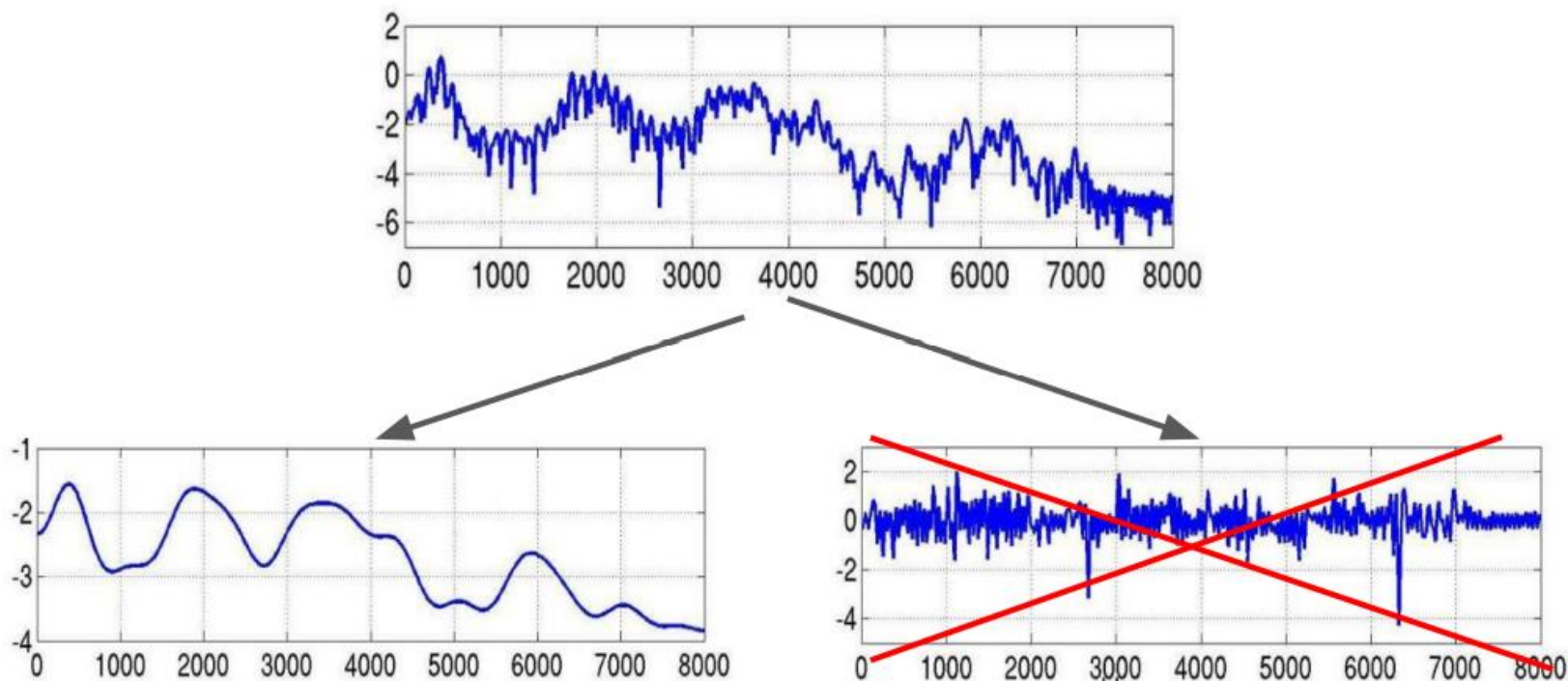
Spectral envelope



Carry identity of
sound



Understanding Cepstrum



An FFT on spectrum referred to as Inverse FFT (IFFT)

MFCC Steps

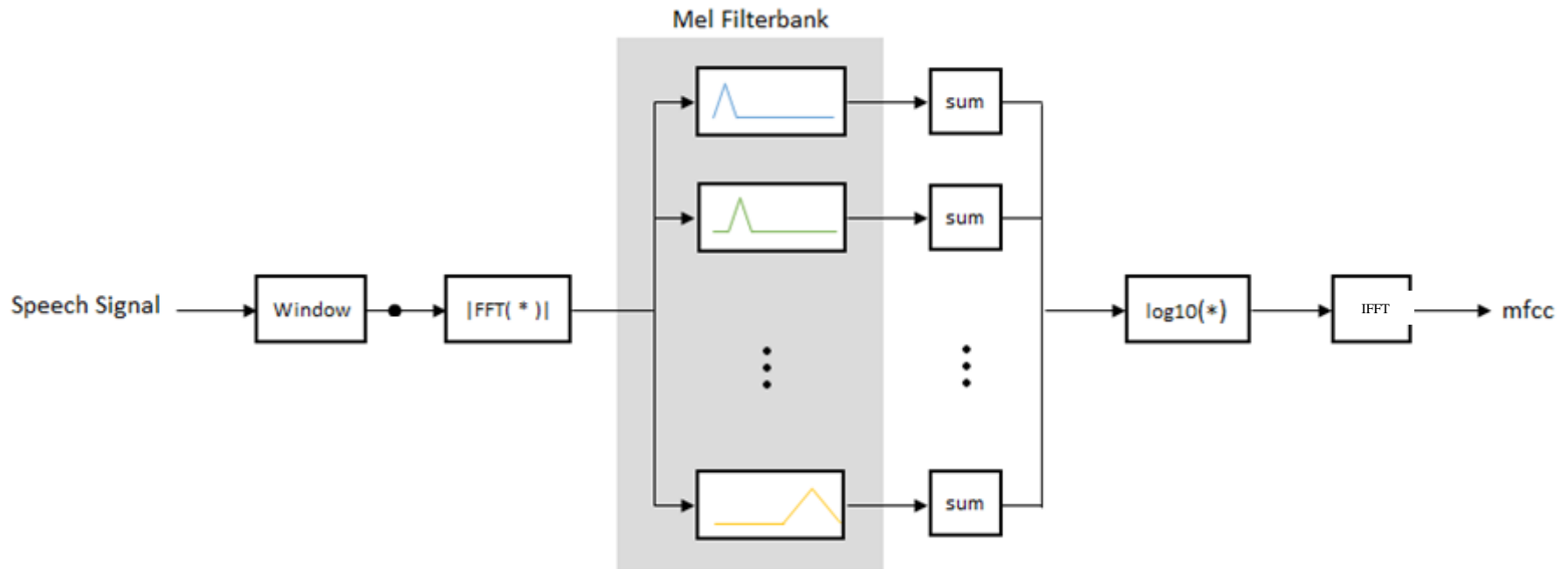
Time-domain signal

$$C(x(t)) = F^{-1} \left[\log \left(F[x(t)] \right) \right]$$

Spectrum

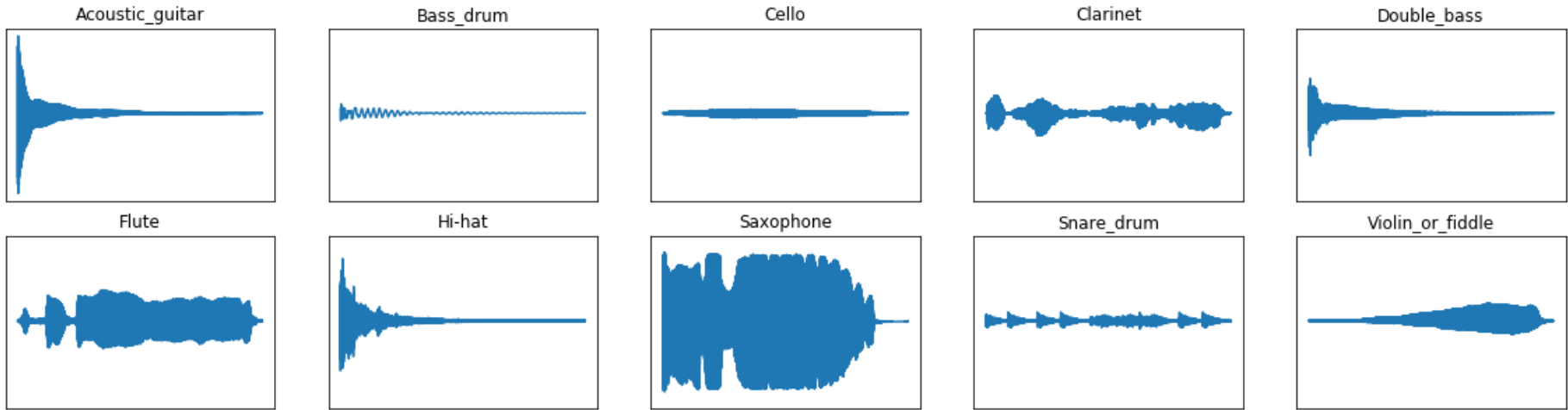
Log spectrum

Cepstrum

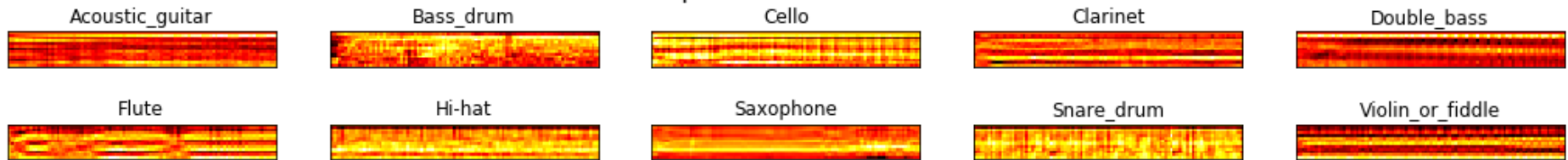


Mel Frequency Cepstrum Coefficients

Time Series

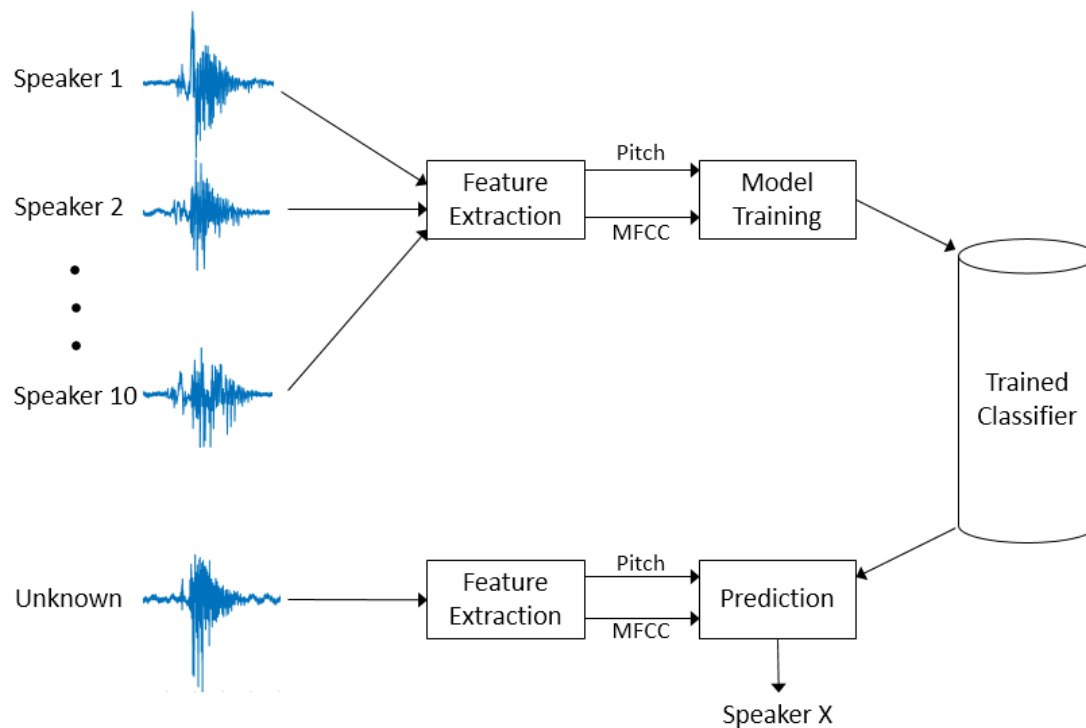
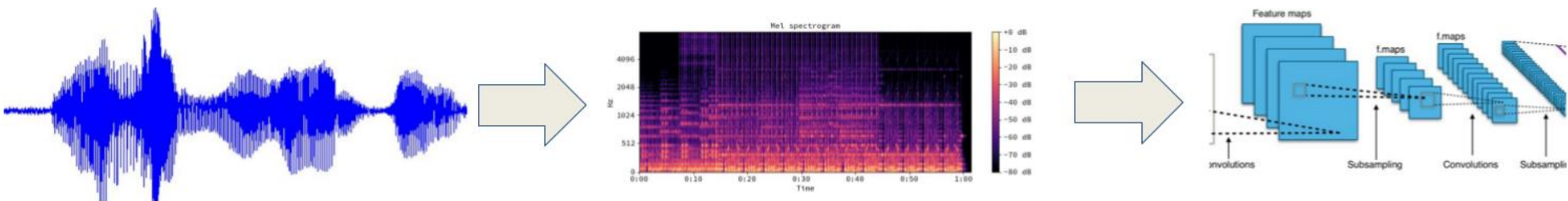


Mel Cepstrum Coefficients

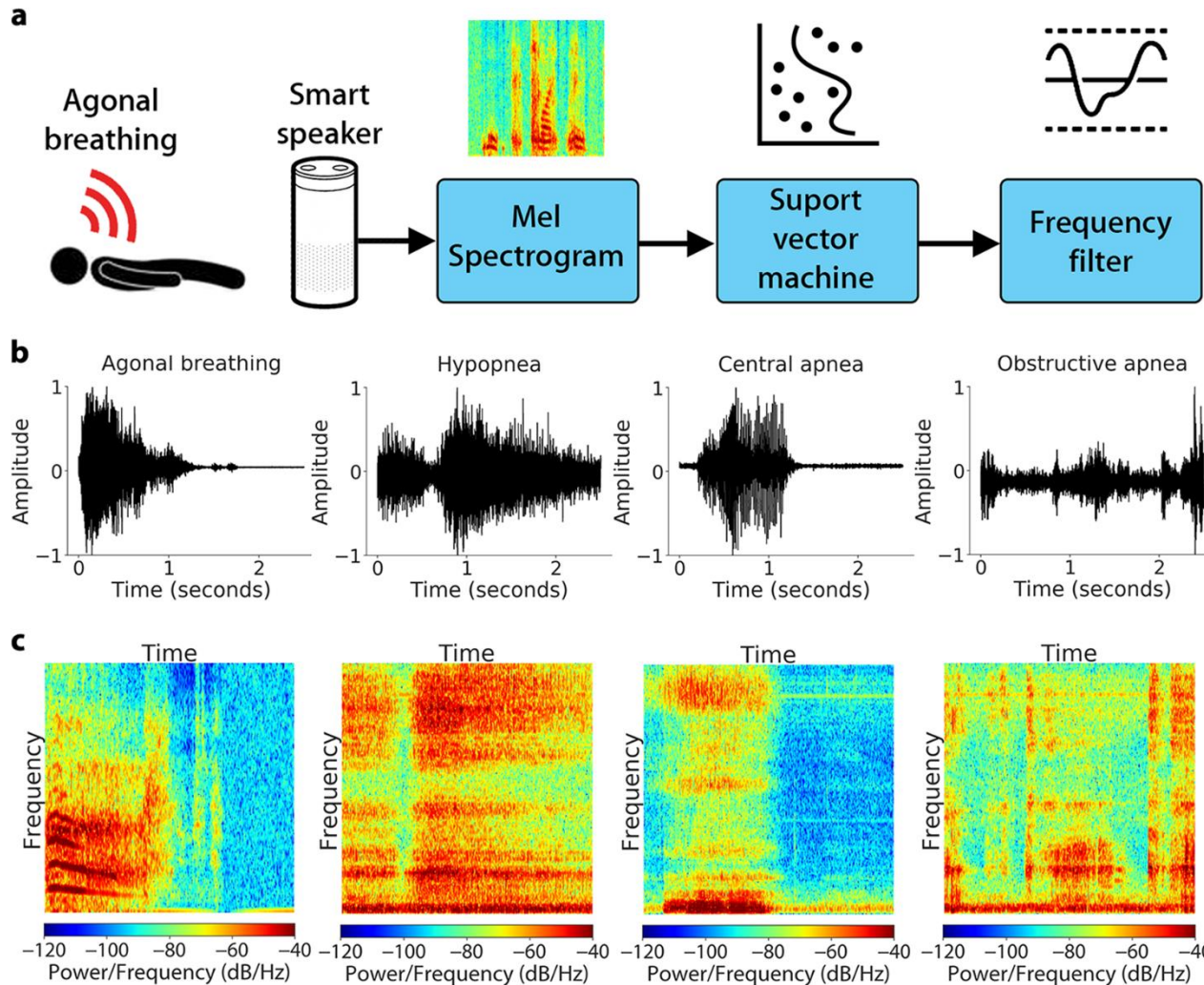


Pretty easy to classify

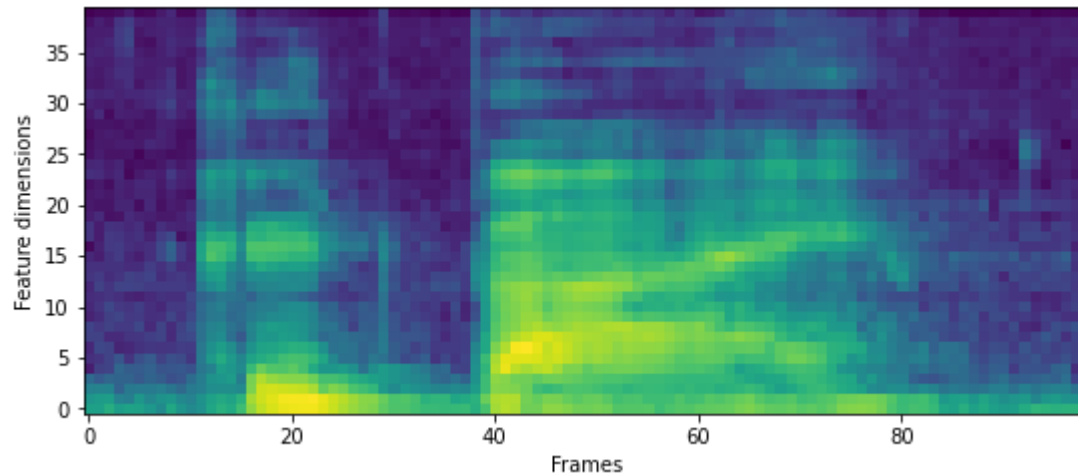
MFCC Use Cases



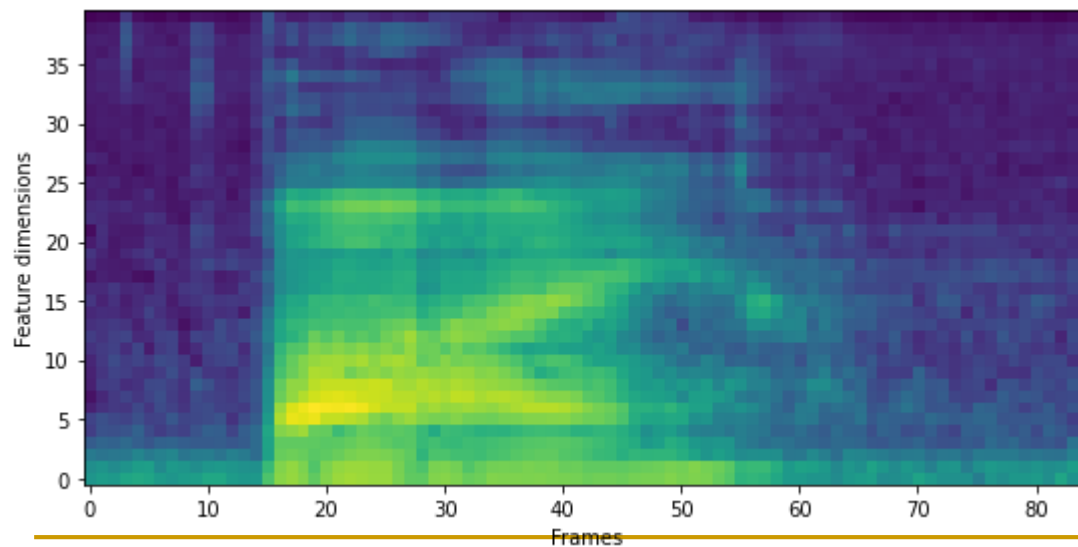
MFCC Use Cases



DTW on MFCC



Goodbye



Bye



DTW on MFCC

