

# Understanding audio data for deep learning

Valerio Velardo

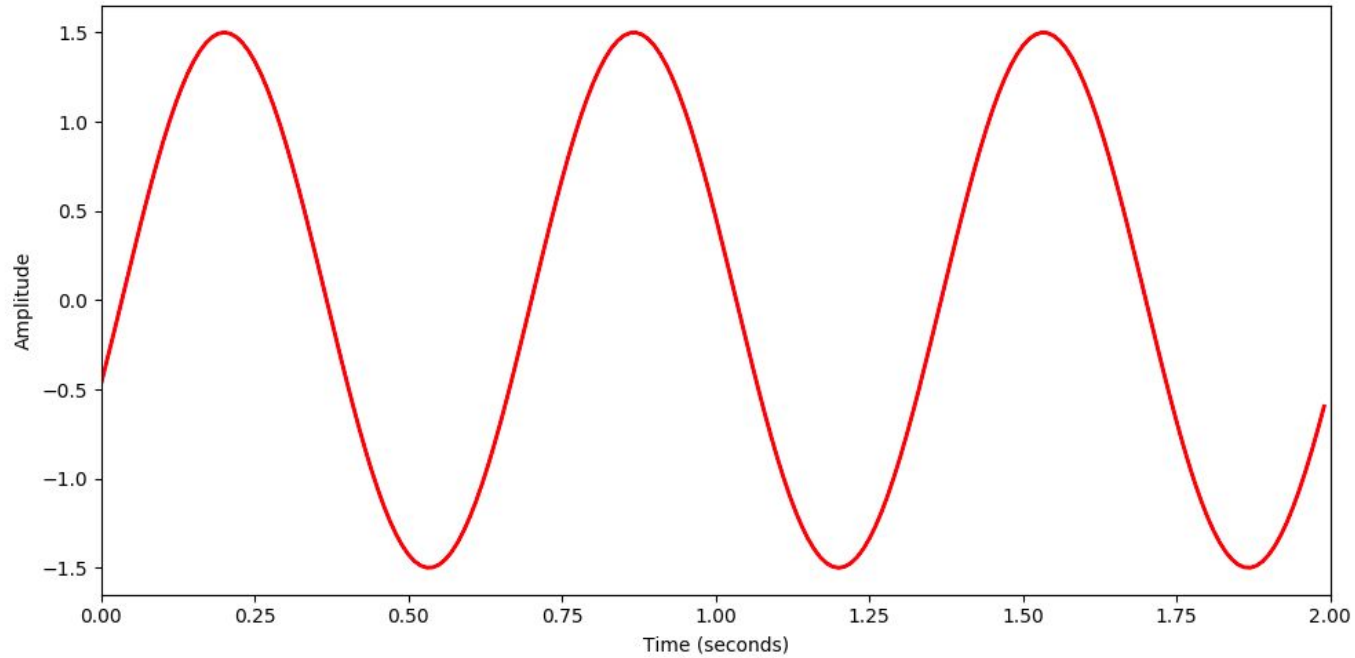
# Sound

---

- Produced by the vibration of an object
- Vibrations determine oscillation of air molecules
- Alternation of air pressure causes a wave

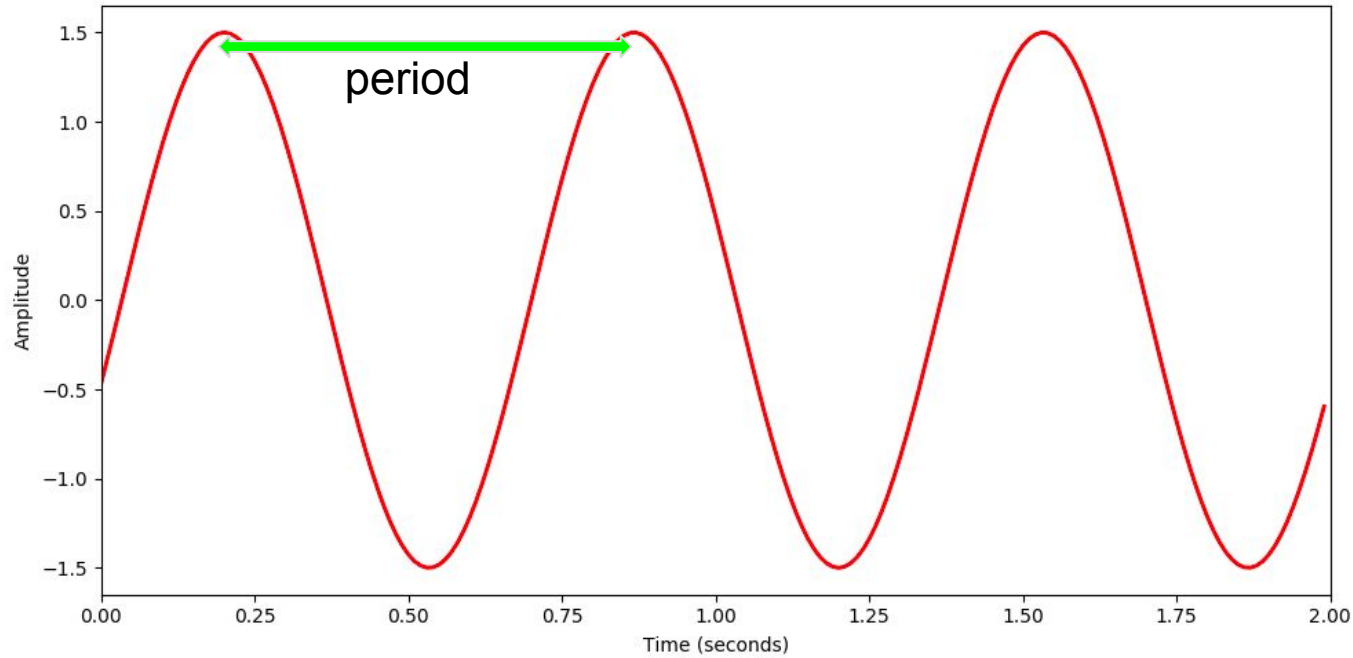
# Waveform

---



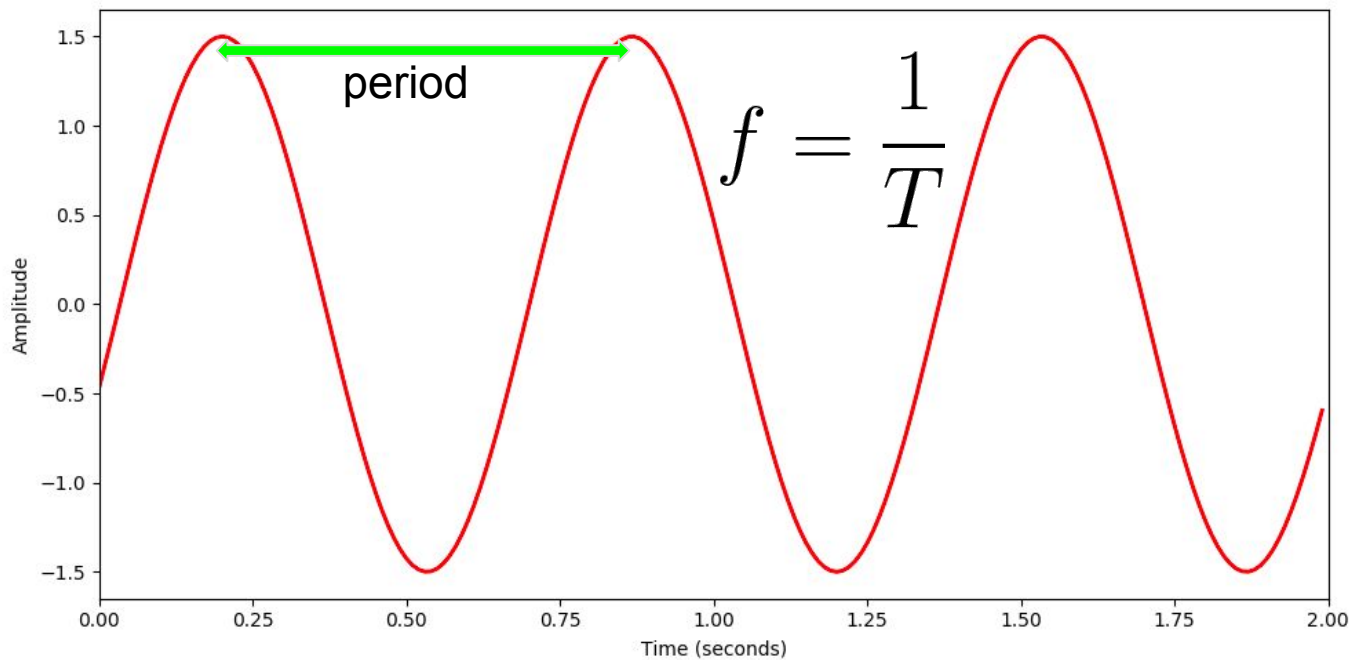
# Waveform

---



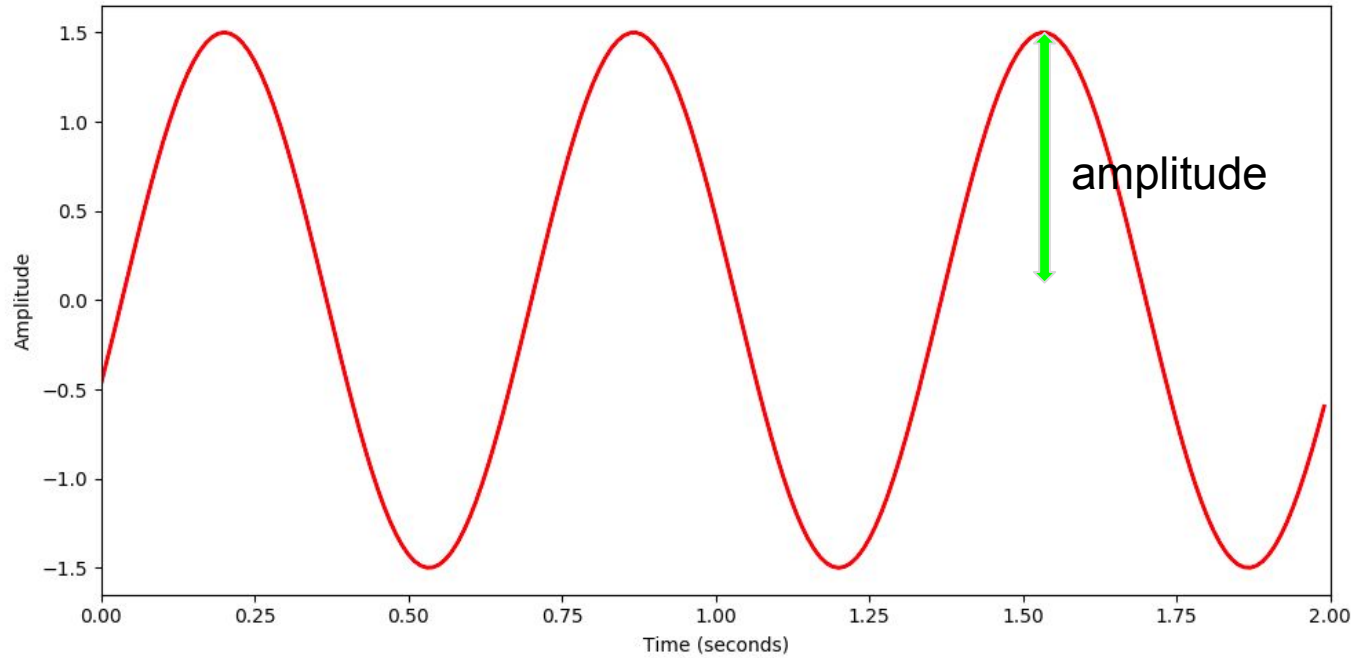
# Waveform

---



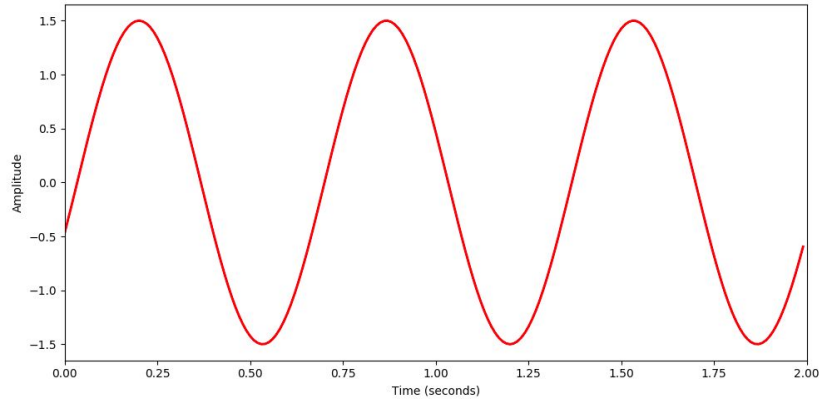
# Waveform

---



# Waveform

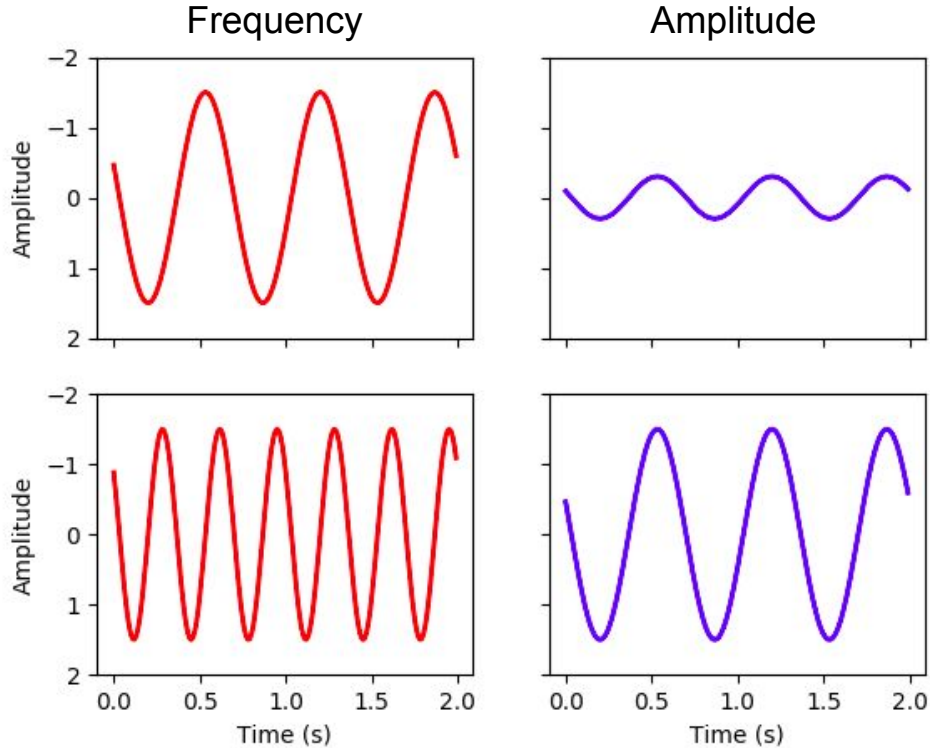
---



$$y(t) = A \sin(2\pi ft + \varphi)$$

# Frequency/pitch and amplitude/loudness

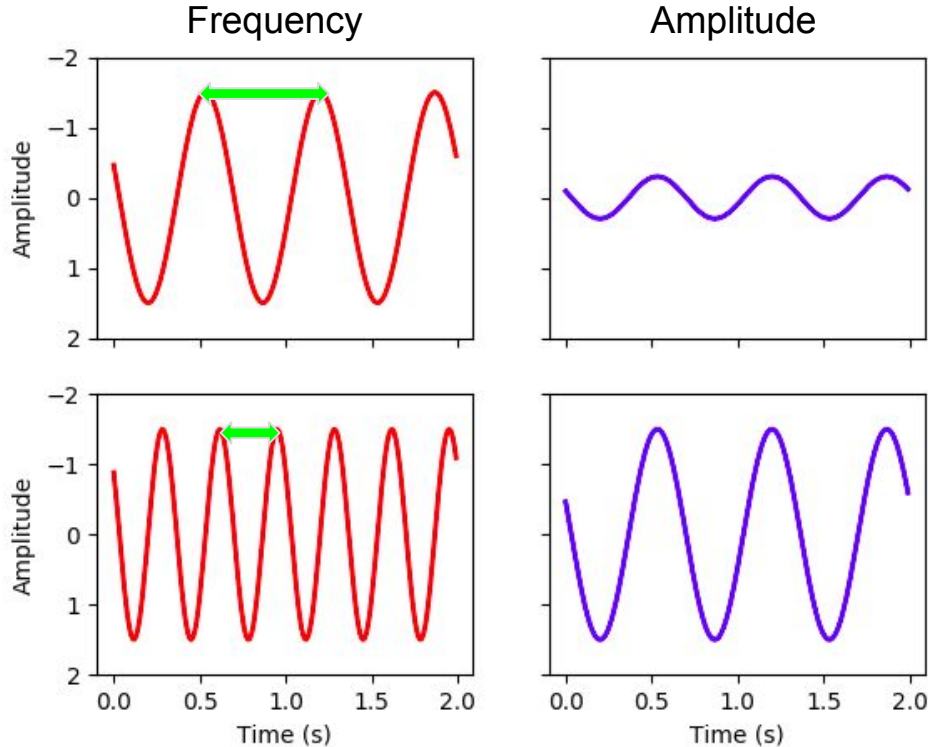
---





# Frequency/pitch and amplitude/loudness

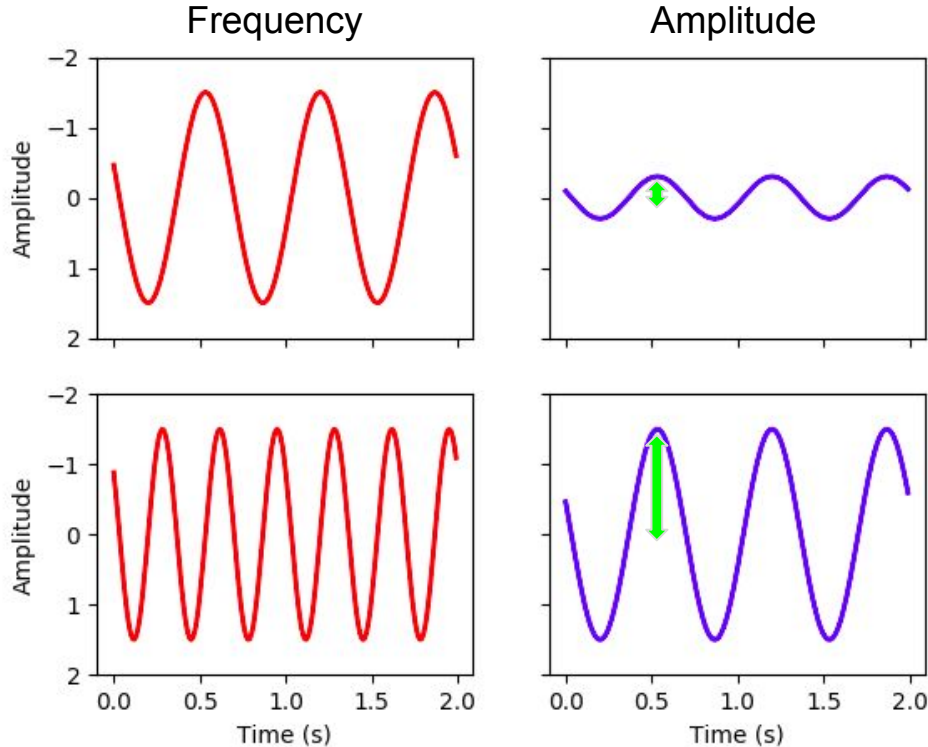
---



higher frequency -> higher pitch

# Frequency/pitch and amplitude/loudness

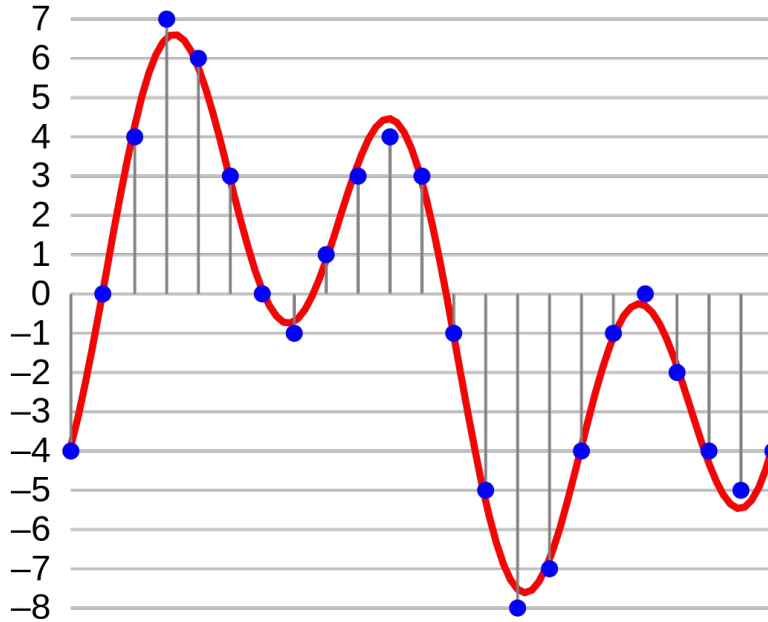
---



larger amplitude -> louder

# Analog digital conversion (ADC)

---



- Signal sampled at uniform time intervals
- Amplitude quantised with limited number of bits

# Analog digital conversion (ADC)

---

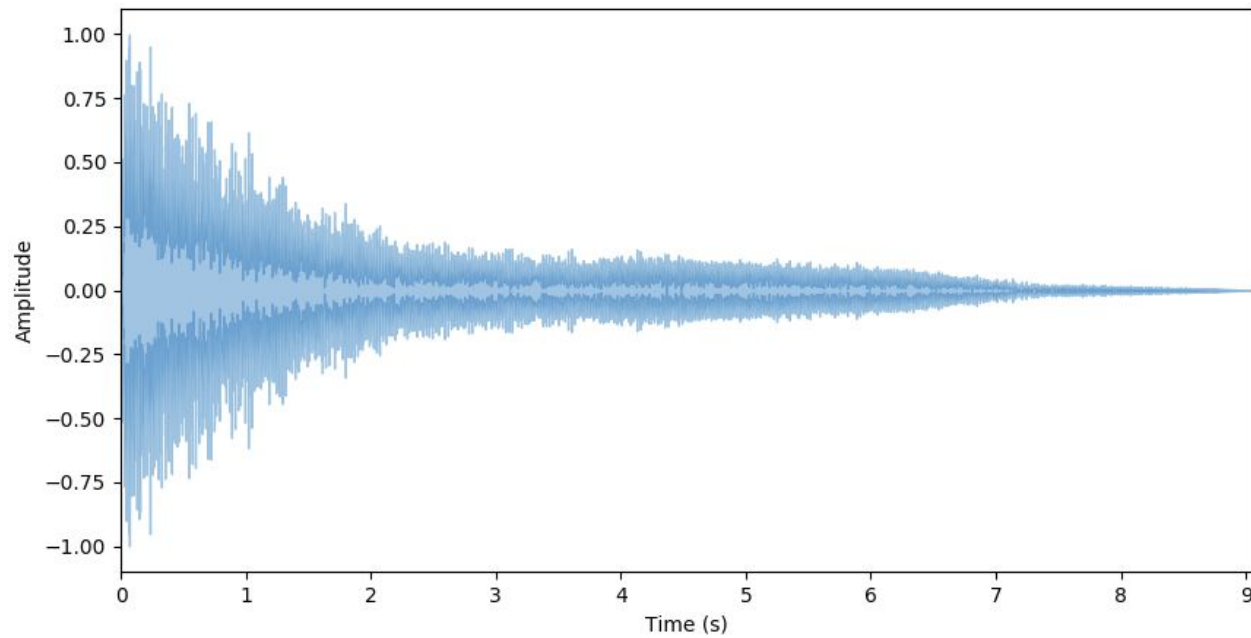


Sample rate = 44,100 Hz

Bit depth = 16 bits/channel

# A real-world sound wave (piano key)

---



# Fourier transform

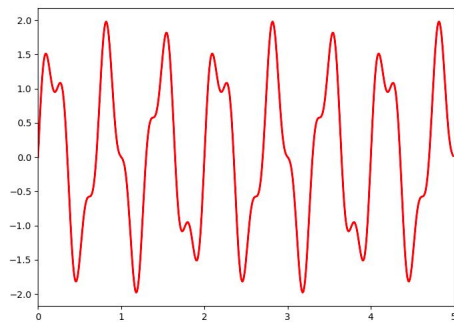
---

- Decompose complex periodic sound into sum of sine waves oscillating at different frequencies

# Fourier transform

---

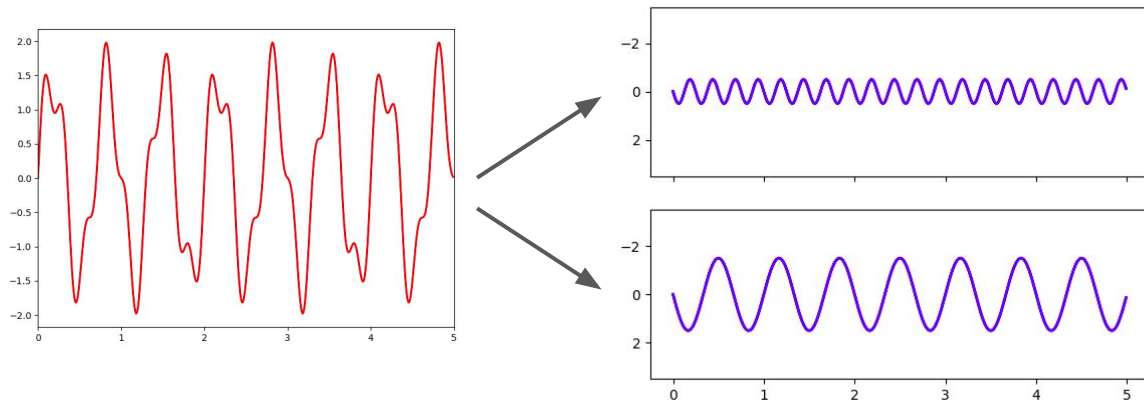
- Decompose complex periodic sound into sum of sine waves oscillating at different frequencies



# Fourier transform

---

- Decompose complex periodic sound into sum of sine waves oscillating at different frequencies

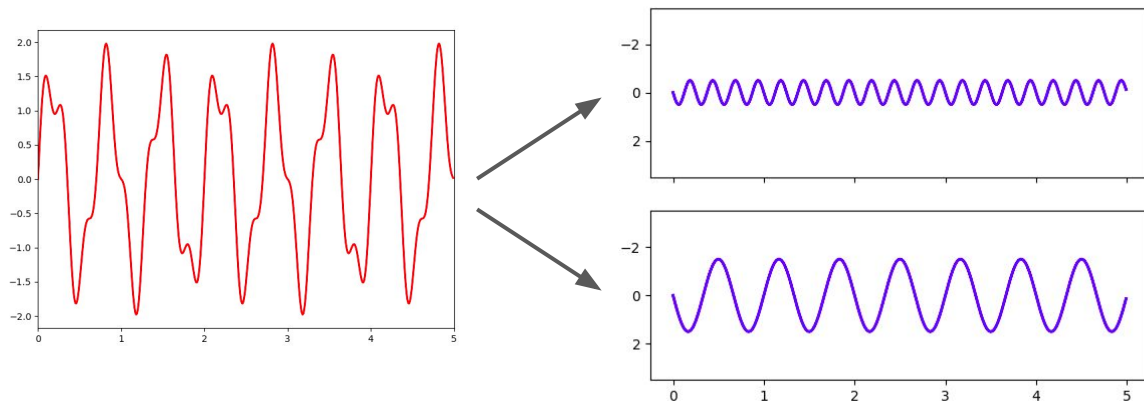




# Fourier transform

---

- Decompose complex periodic sound into sum of sine waves oscillating at different frequencies

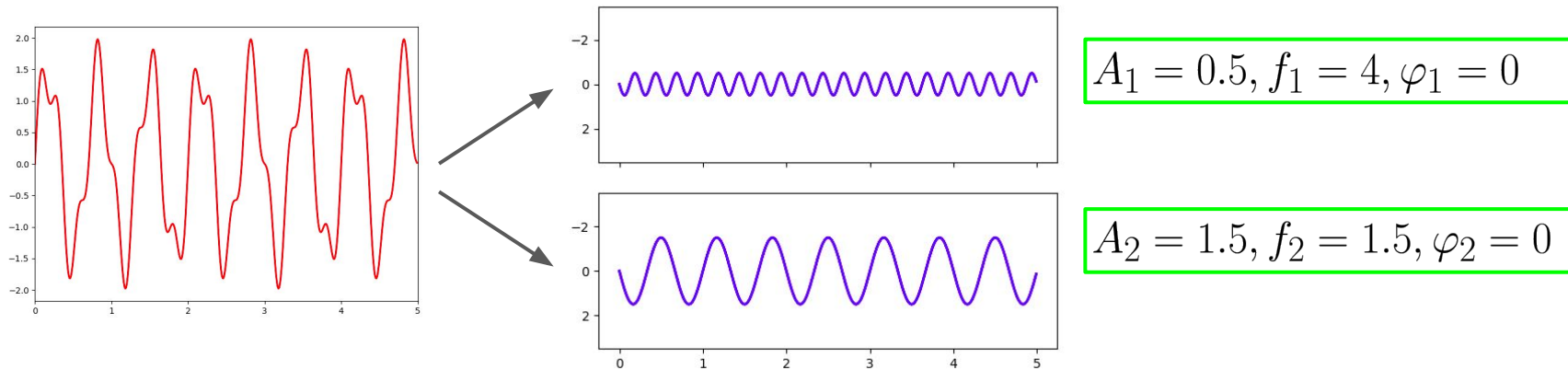


$$s = A_1 \sin(2\pi f_1 t + \varphi_1) + A_2 \sin(2\pi f_2 t + \varphi_2)$$

# Fourier transform

---

- Decompose complex periodic sound into sum of sine waves oscillating at different frequencies

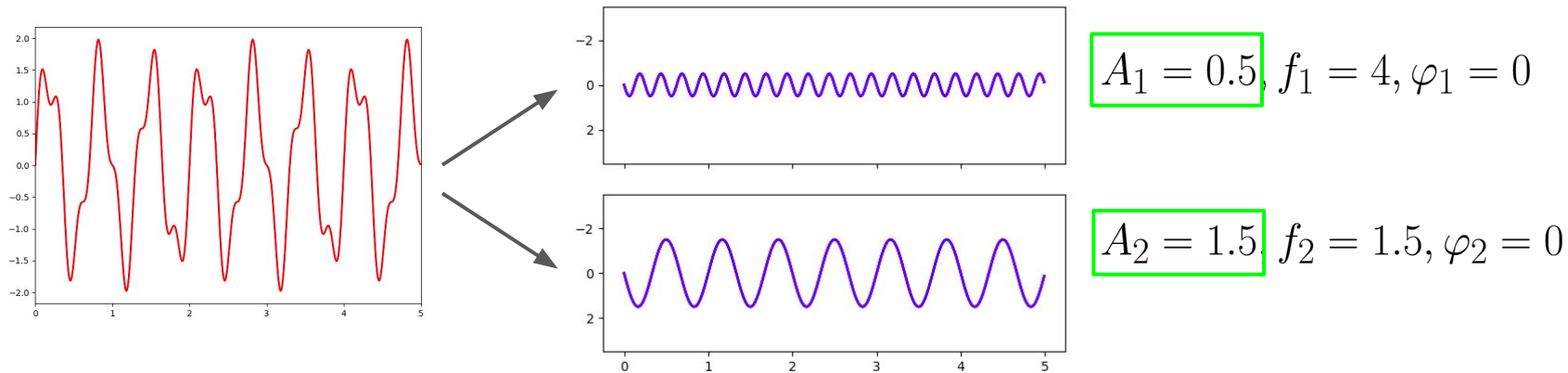


$$s = A_1 \sin(2\pi f_1 t + \varphi_1) + A_2 \sin(2\pi f_2 t + \varphi_2)$$

# Fourier transform

---

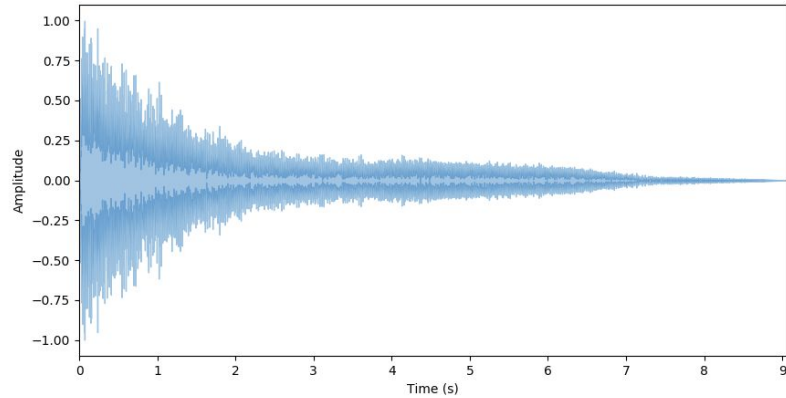
- Decompose complex periodic sound into sum of sine waves oscillating at different frequencies



$$s = A_1 \sin(2\pi f_1 t + \varphi_1) + A_2 \sin(2\pi f_2 t + \varphi_2)$$

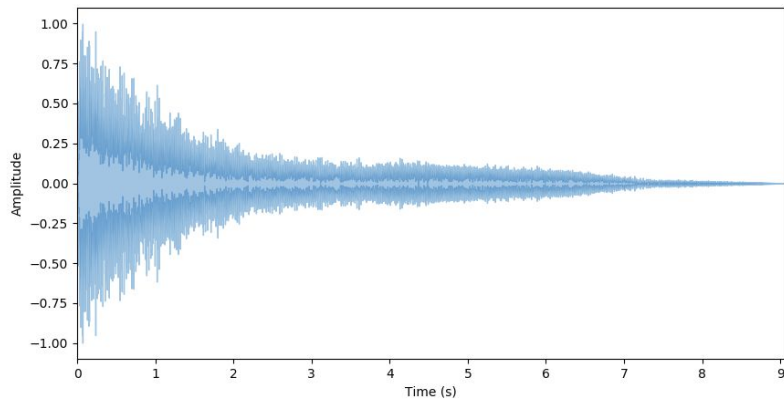
# Fourier transform

---

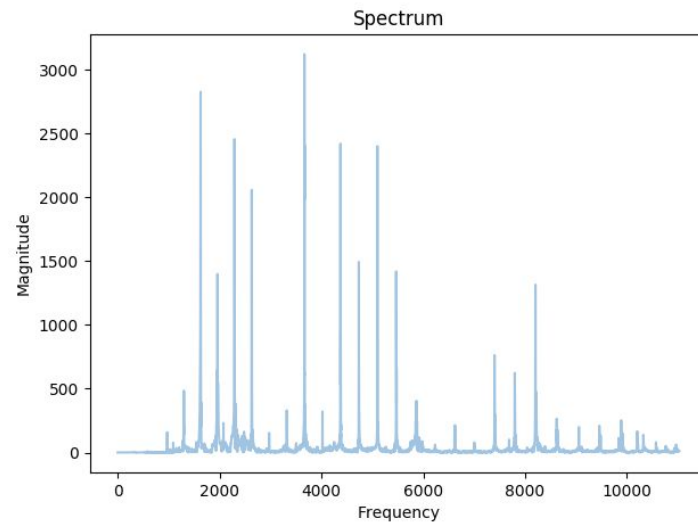


# Fourier transform

---

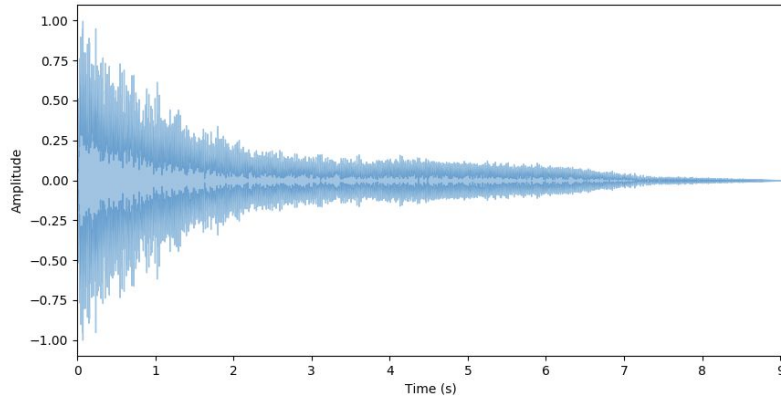


FFT  
→

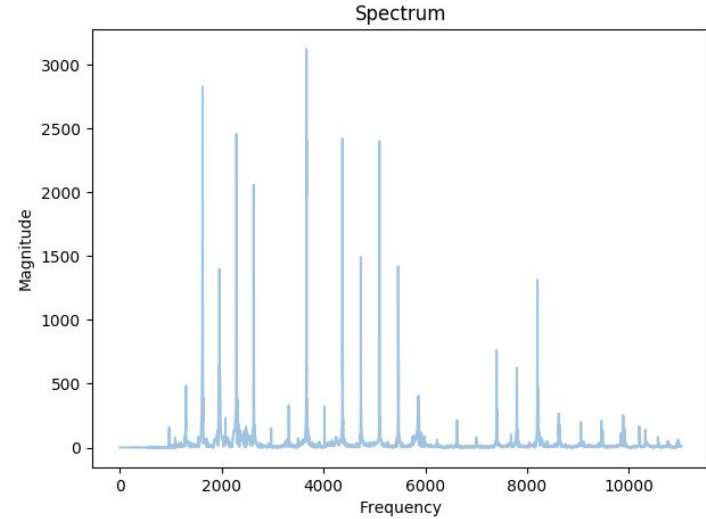


# Fourier transform

---



FFT  
→



- From *time domain* to *frequency domain*
- No time information

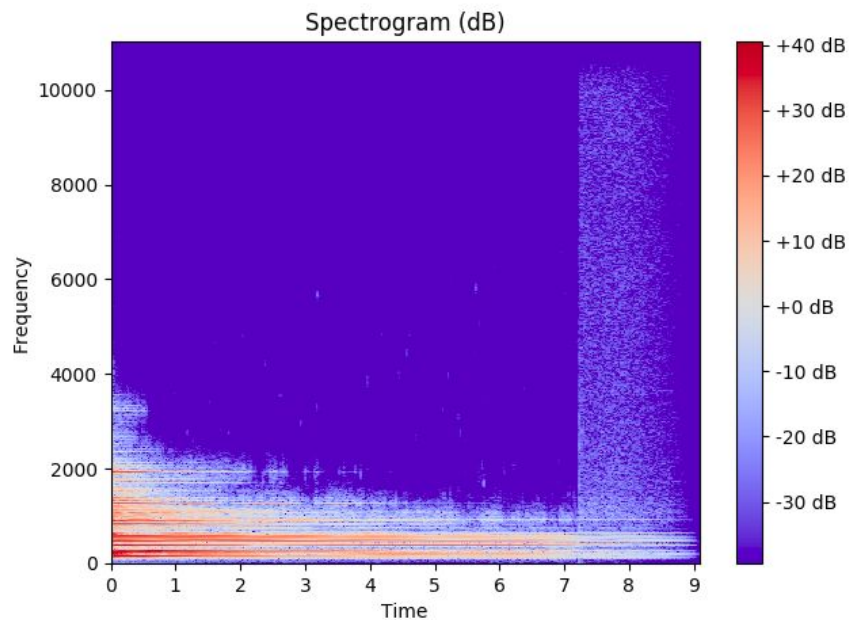
# Short Time Fourier Transform (STFT)

---

- Computes several FFT at different intervals
- Preserves time information
- Fixed frame size (e.g., 2048 samples)
- Gives a *spectrogram* (time + frequency + magnitude)

# Short Time Fourier Transform (STFT)

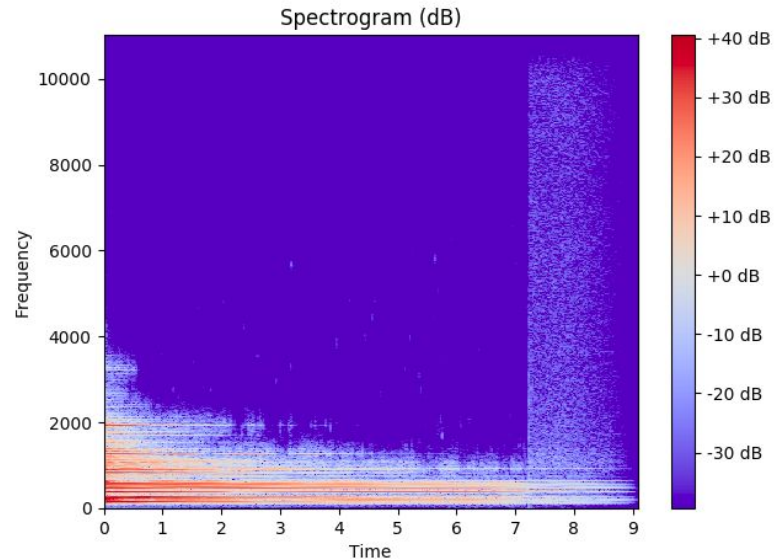
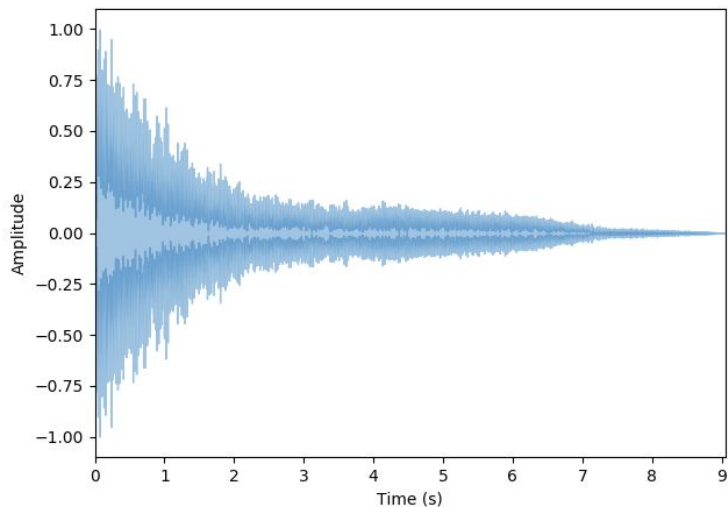
---





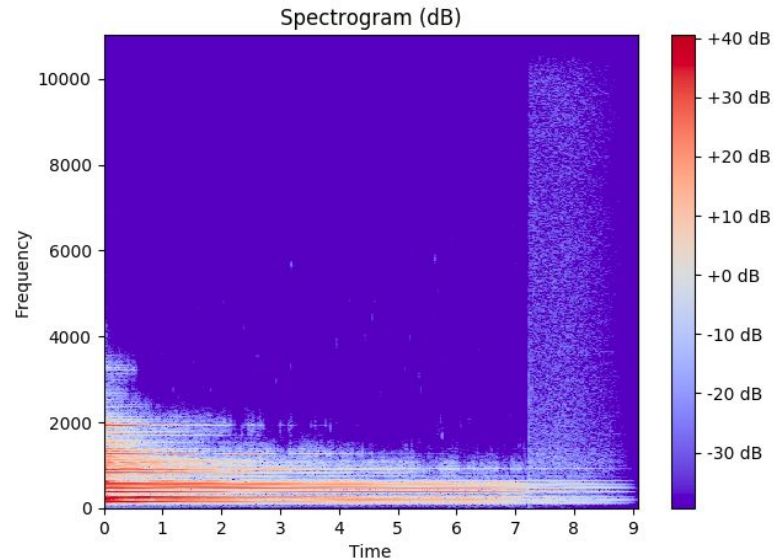
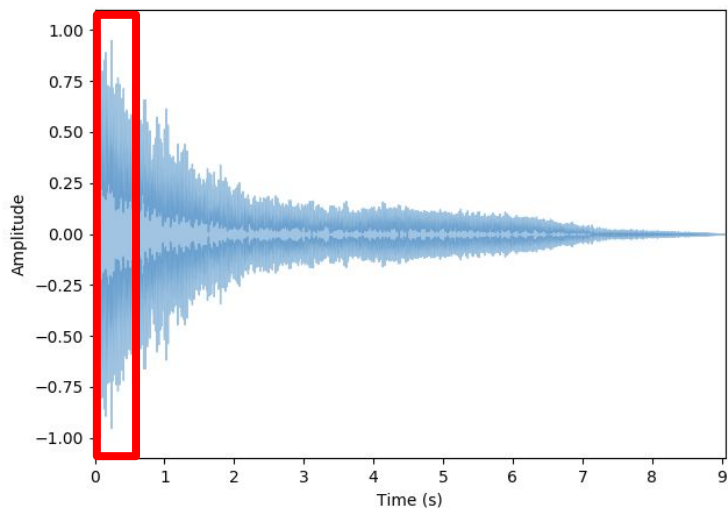
# Short Time Fourier Transform (STFT)

---

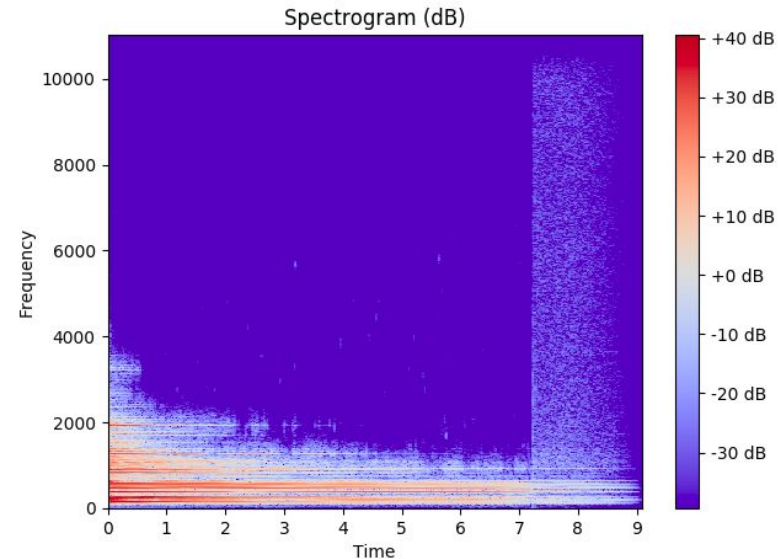
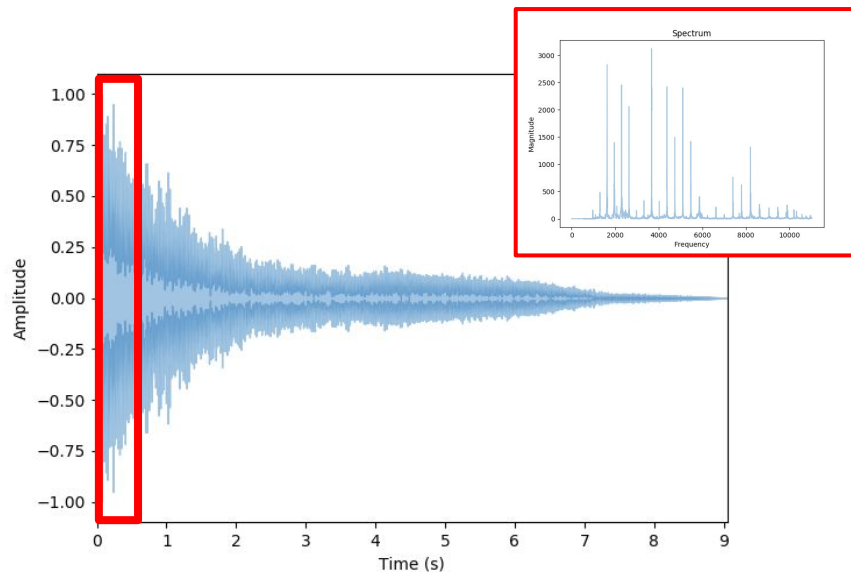


# Short Time Fourier Transform (STFT)

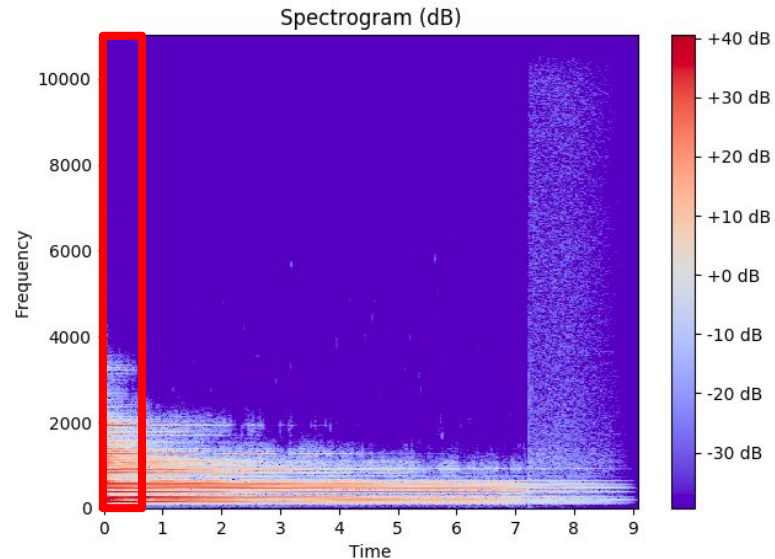
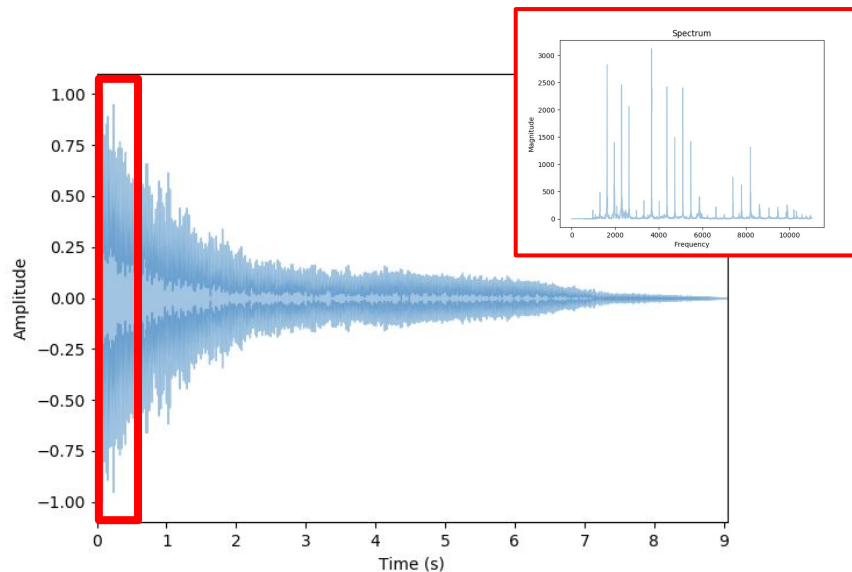
---



# Short Time Fourier Transform (STFT)

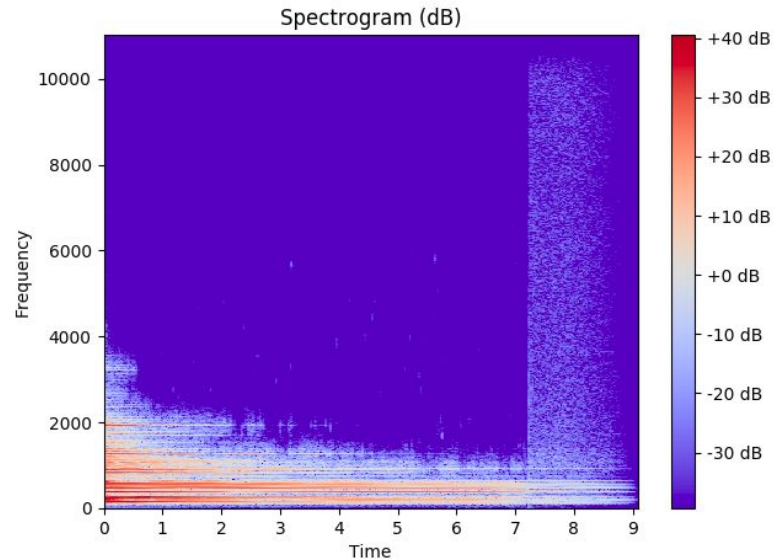
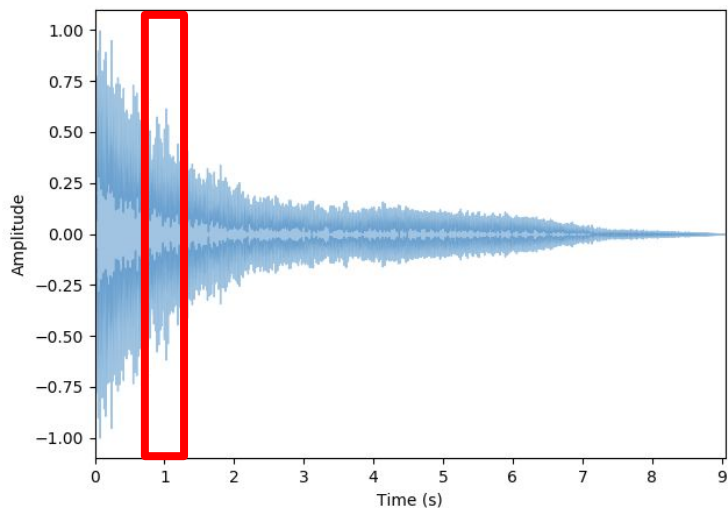


# Short Time Fourier Transform (STFT)



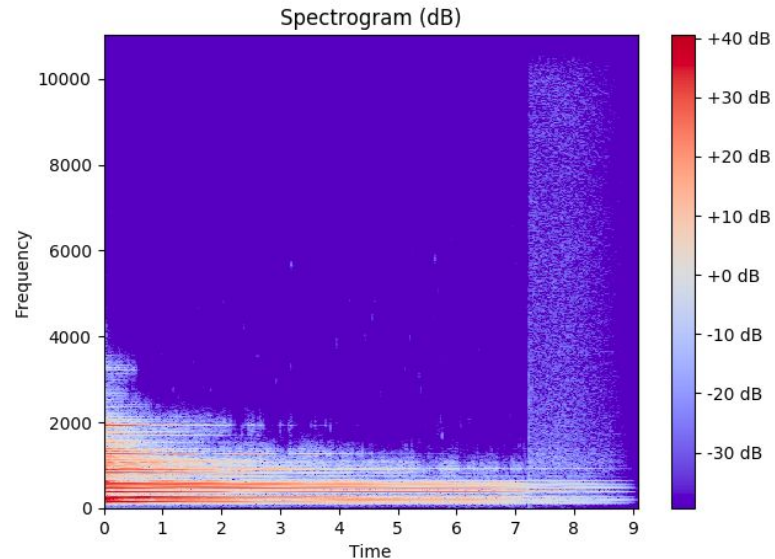
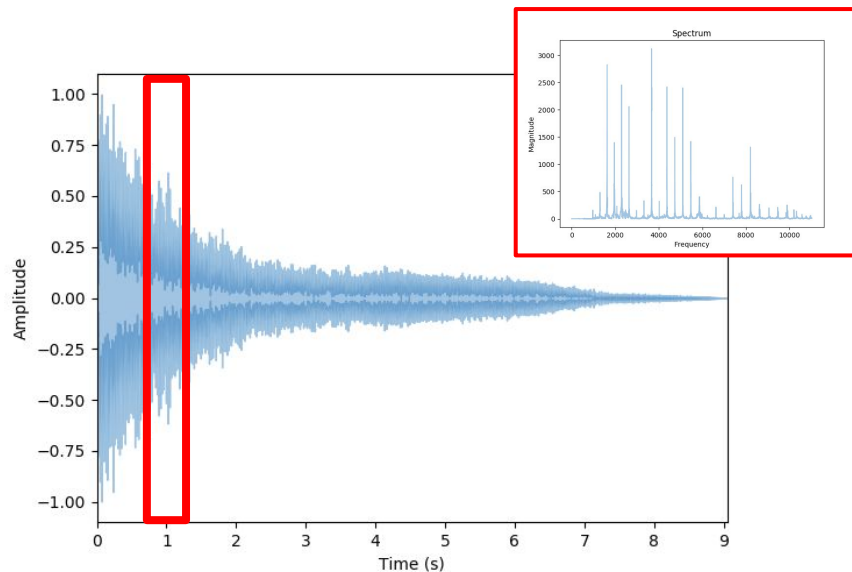
# Short Time Fourier Transform (STFT)

---



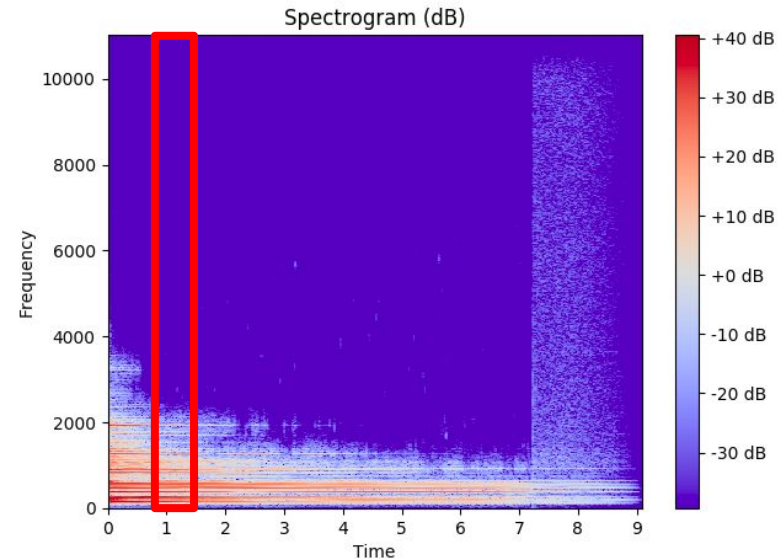
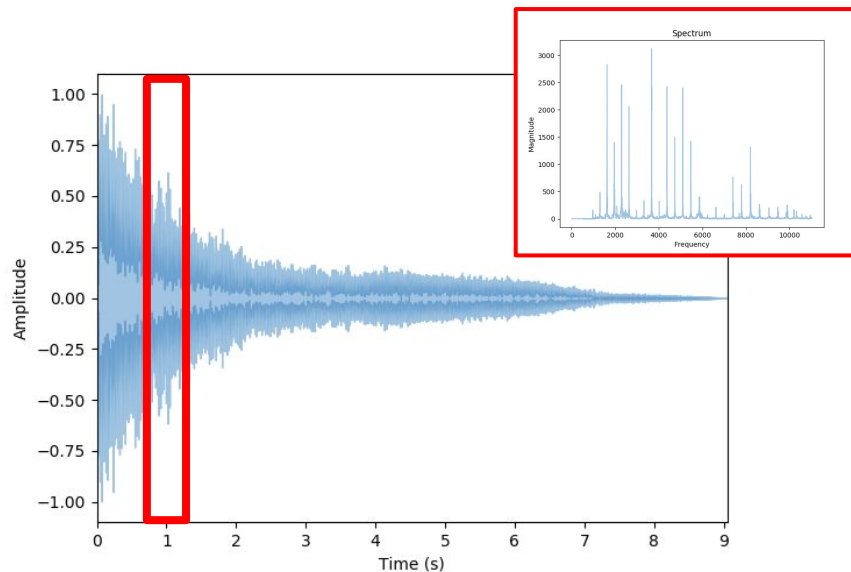
# Short Time Fourier Transform (STFT)

---

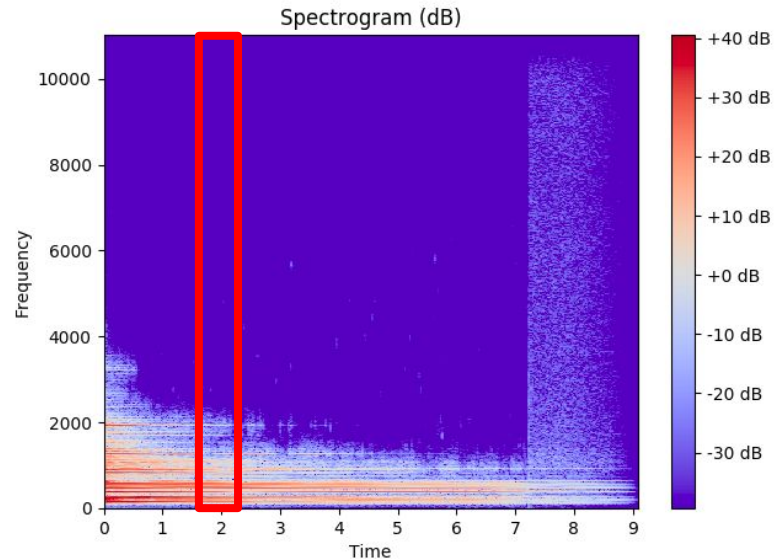
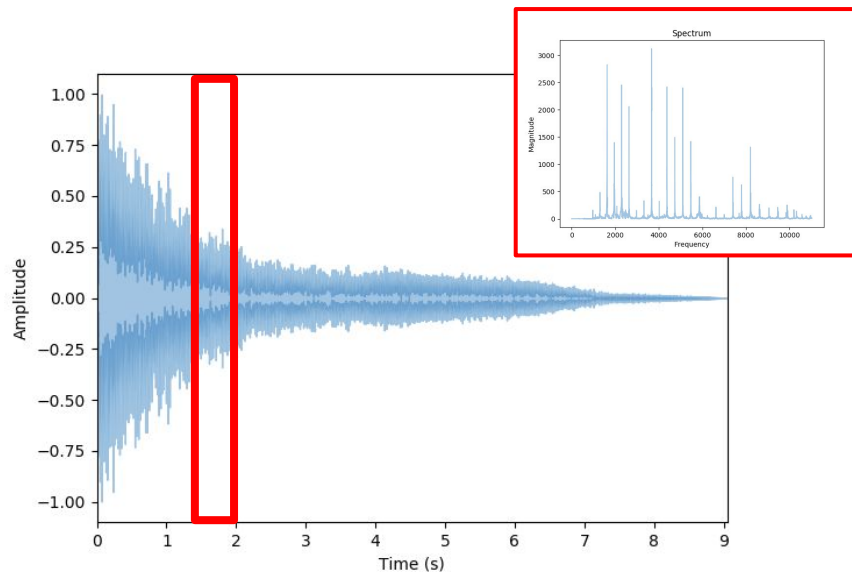


# Short Time Fourier Transform (STFT)

---



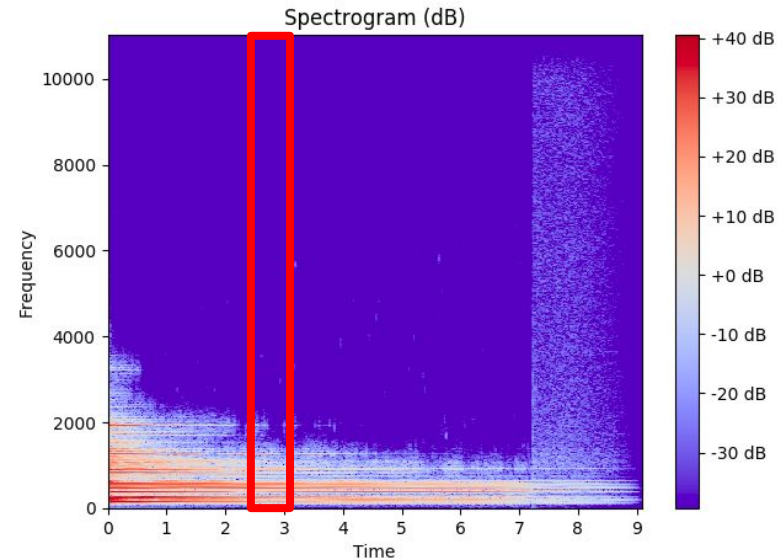
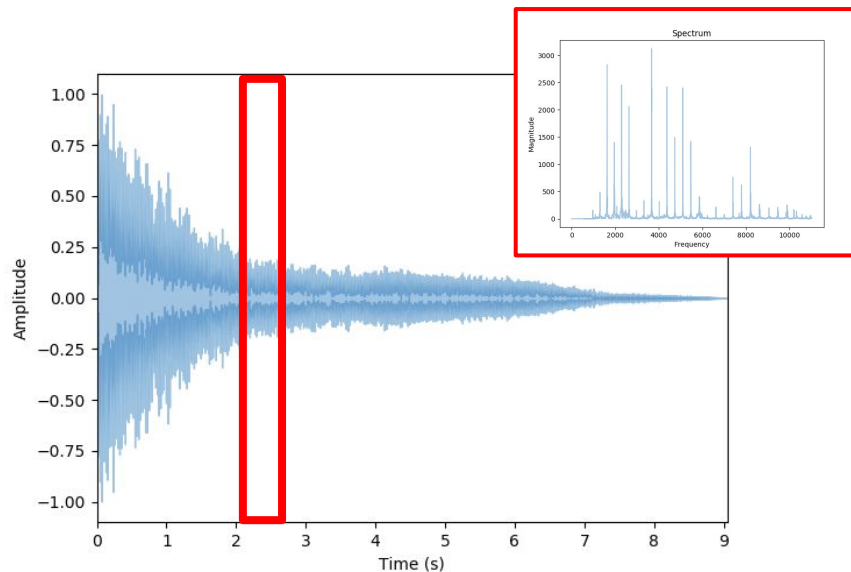
# Short Time Fourier Transform (STFT)





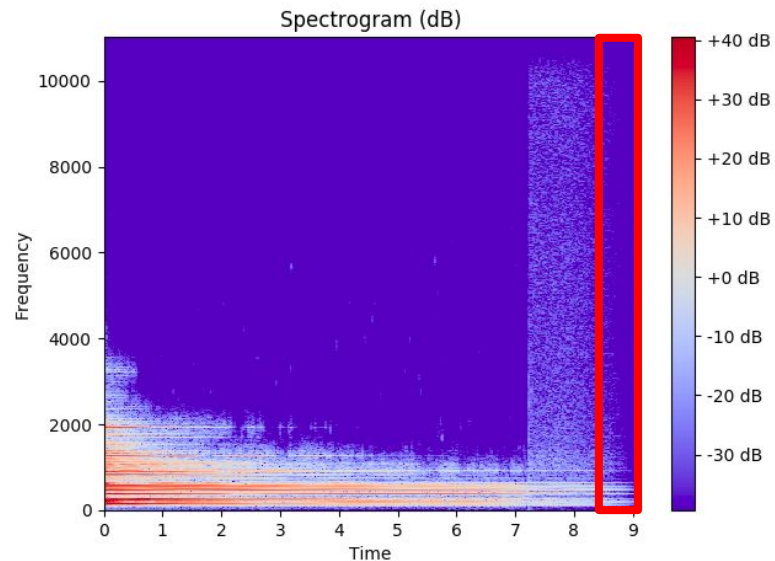
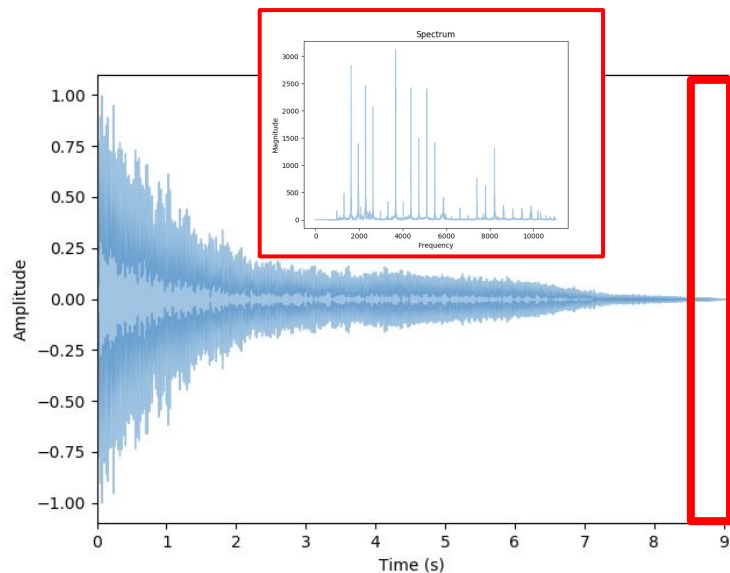
# Short Time Fourier Transform (STFT)

---



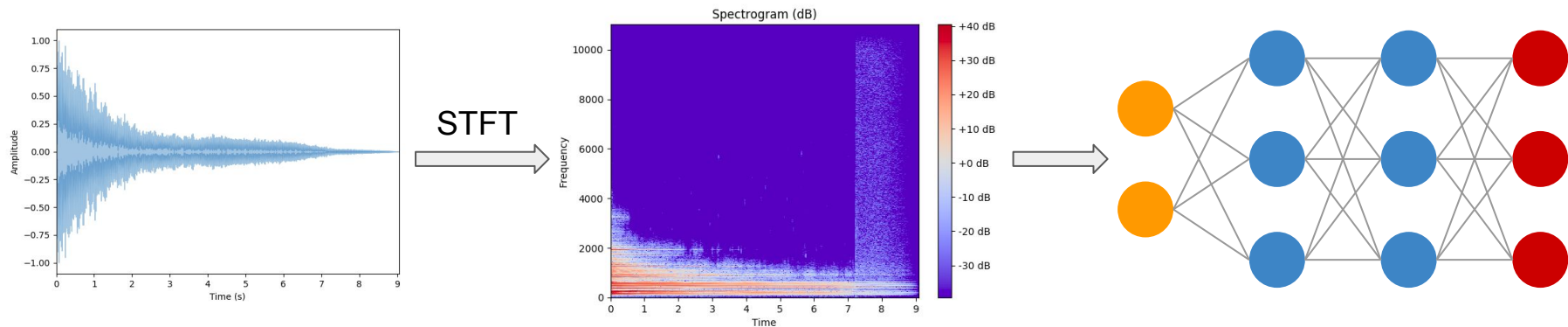
# Short Time Fourier Transform (STFT)

---



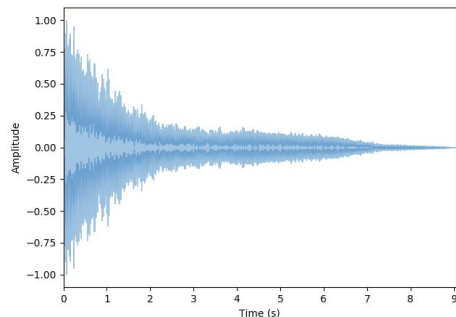
# DL pre-processing pipeline for audio data

---

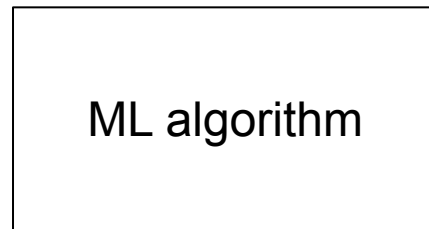


# Traditional ML pre-processing pipeline for audio data

---



- Amplitude envelope
- Spectral centroid
- 0-crossing rate
- Spectral flux
- ...



- Feature engineering
- Perform STFT
- Extract time + frequency domain features

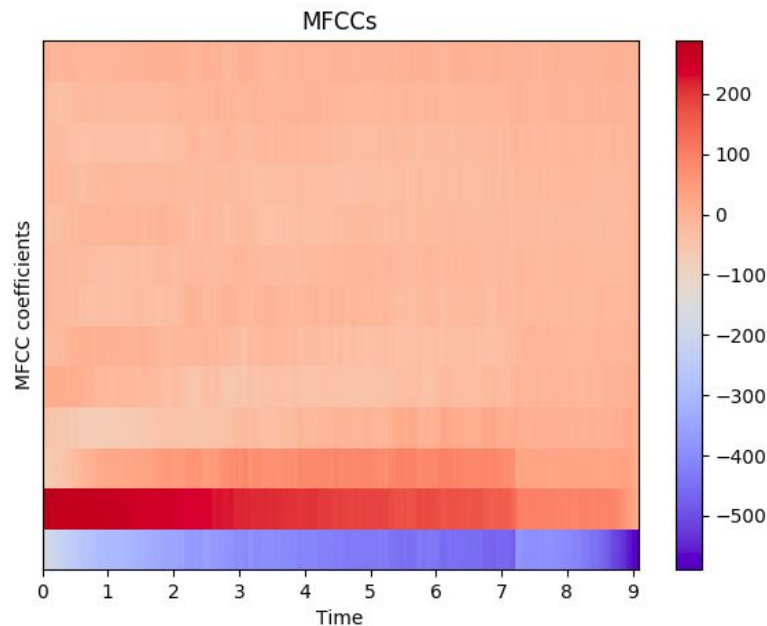
# Mel Frequency Cepstral Coefficients (MFCCs)

---

- Capture timbral/textural aspects of sound
- Frequency domain feature
- Approximate human auditory system
- 13 to 40 coefficients
- Calculated at each frame

# Mel Frequency Cepstral Coefficients (MFCCs)

---



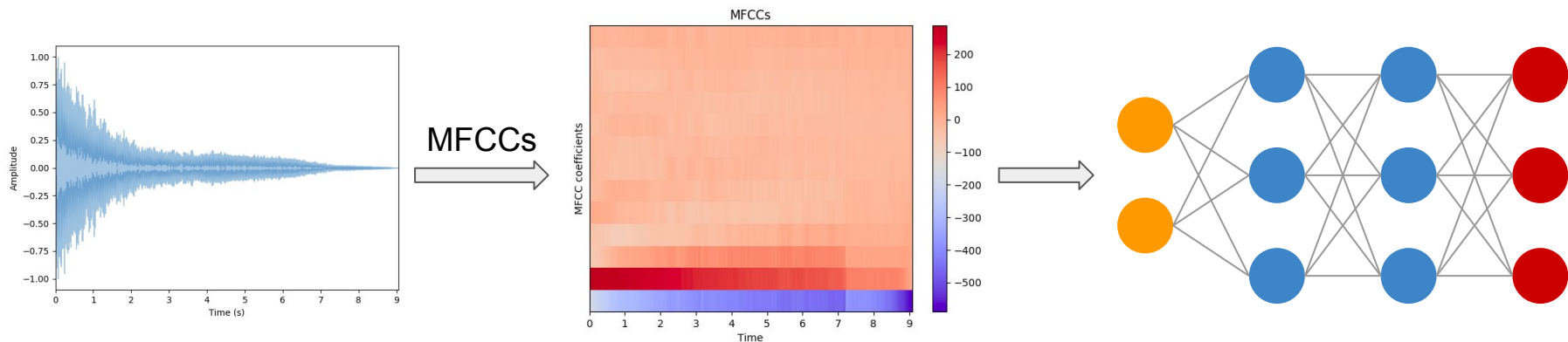
# MFCCs applications

---

- Speech recognition
- Music genre classification
- Music instrument classification
- ...

# DL pre-processing pipeline for audio data

---





# What's up next?

---

- Perform FFT and STFT with Python
- Extract MFCCs
- Get familiar with *librosa*