# Mel-Frequency Cepstral Coefficients Explained Easily
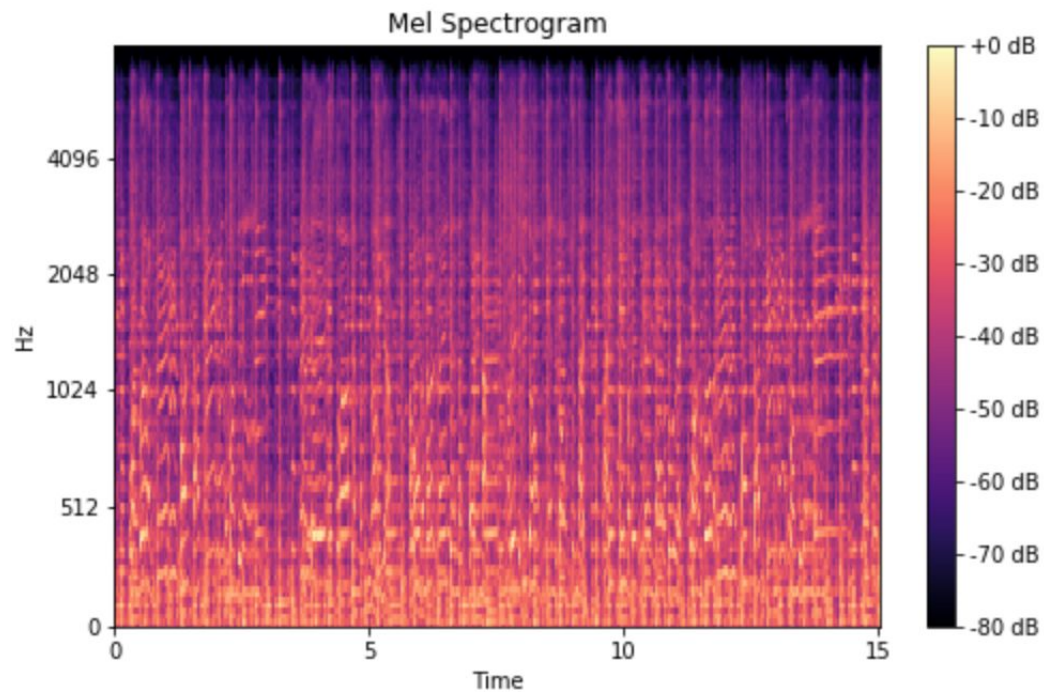
Valerio Velardo

# Join the community!



thesoundofai.slack.com

# Previously...



Mel Spectrogram

# Mel-Frequency Cepstral Coefficients

# Mel-Frequency Cepstral Coefficients

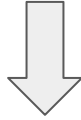# Mel-Frequency Cepstral Coefficients

Mel-Frequency <span style="color:red">Cepstral</span> Coefficients

# Cepstrum

Cepstrum

# Cepstrum

⬇

# Spectrum

Cepstrum

⇩

Spectrum

Cepstrum

⇓

Spectrum

Quefrency    Liftering    Rhamonic

| Cepstrum | Quefrency | Liftering | Rhamonic |
| --- | --- | --- | --- |
| ⇩ | ⇩ | ⇩ | ⇩ |
| Spectrum | Frequency | Filtering | Harmonic |

# An historical note on Cepstrum

- Developed while studying echoes in seismic signals (1960s)

- Audio feature of choice for speech recognition / identification (1970s)

- Music processing (2000s)

# Computing the cepstrum

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

# Computing the cepstrum

Time-domain signal

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

# Computing the cepstrum

Time-domain signal

Spectrum

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

# Computing the cepstrum

Time-domain signal

Spectrum

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

Log spectrum

# Computing the cepstrum

Time-domain signal

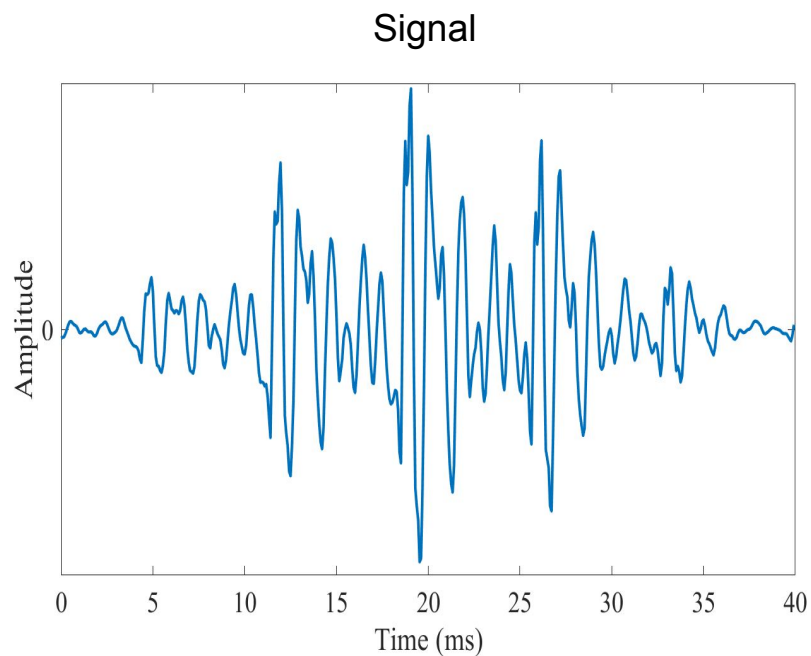$$C(x(t)) = F^{-1}[log(F[x(t)])]$$

Spectrum

Log spectrum

Cepstrum

Spectrum
of
a spectrum

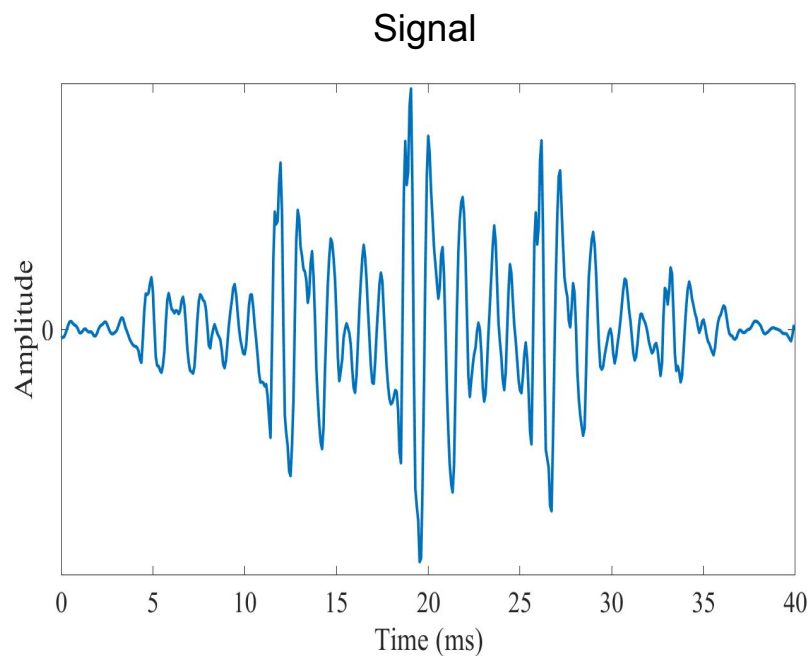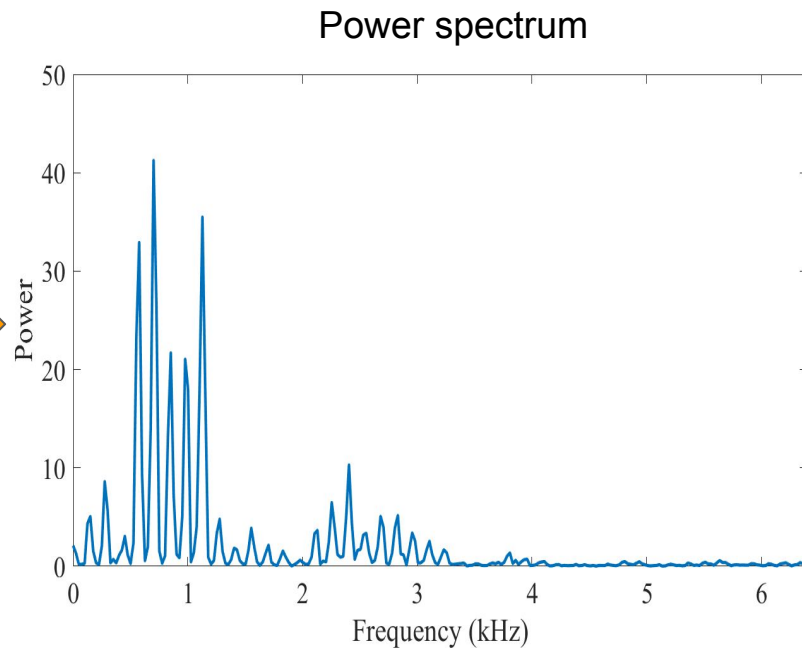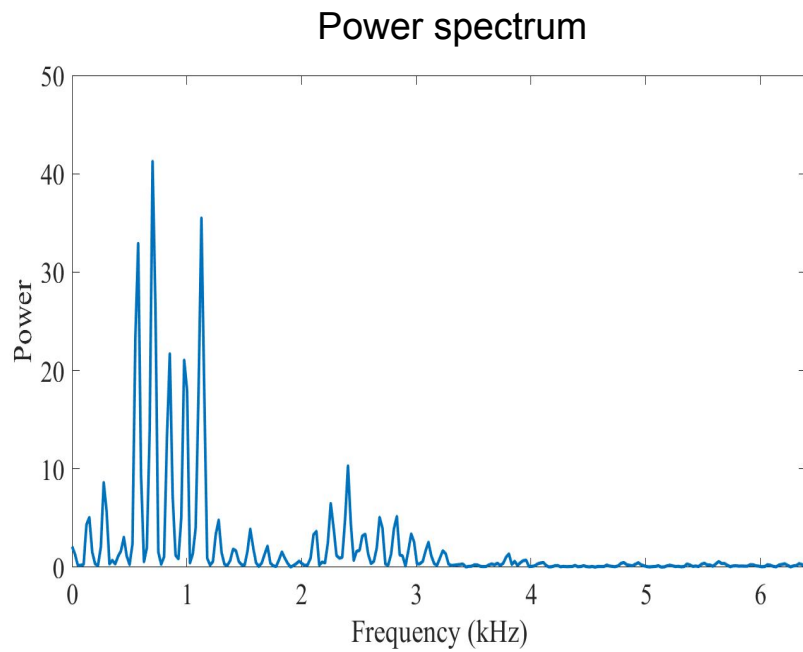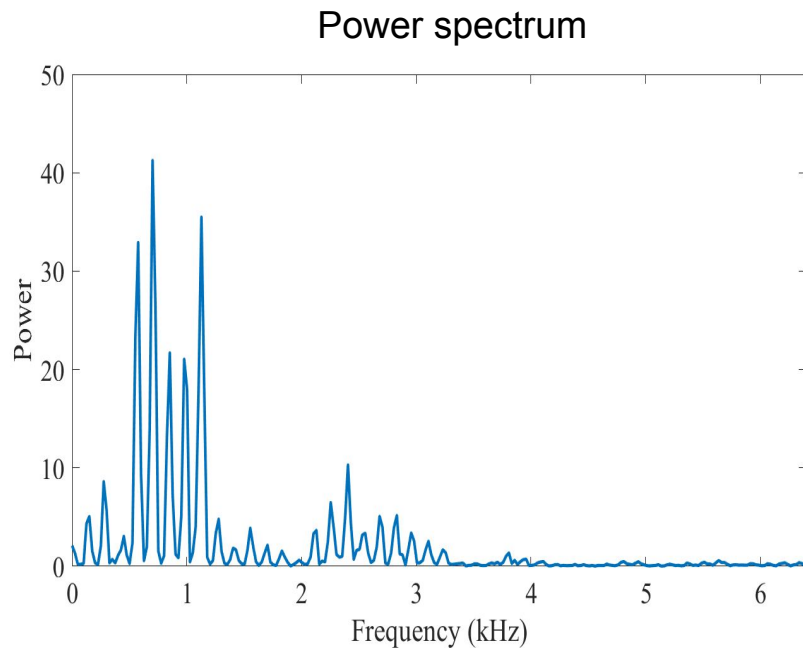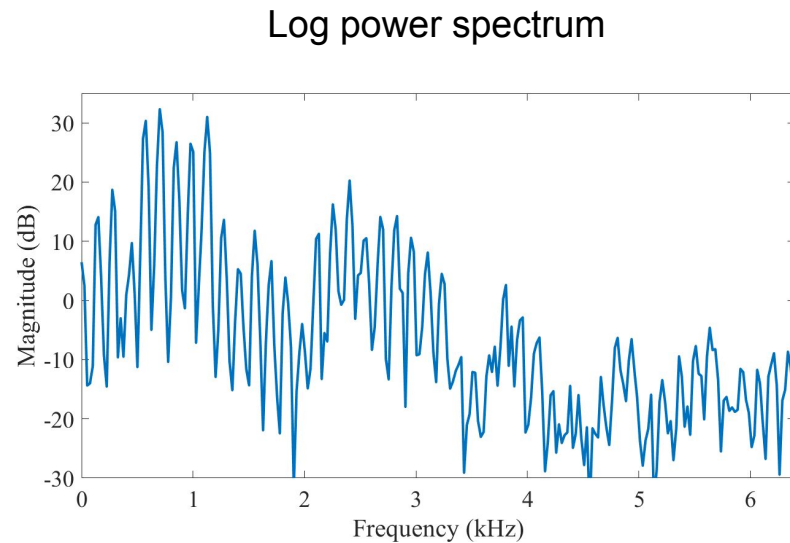Cepstrum

# Visualising the cepstrum



Signal

# Visualising the cepstrum



Signal

DFT

Power spectrum

# Visualising the cepstrum



Power spectrum

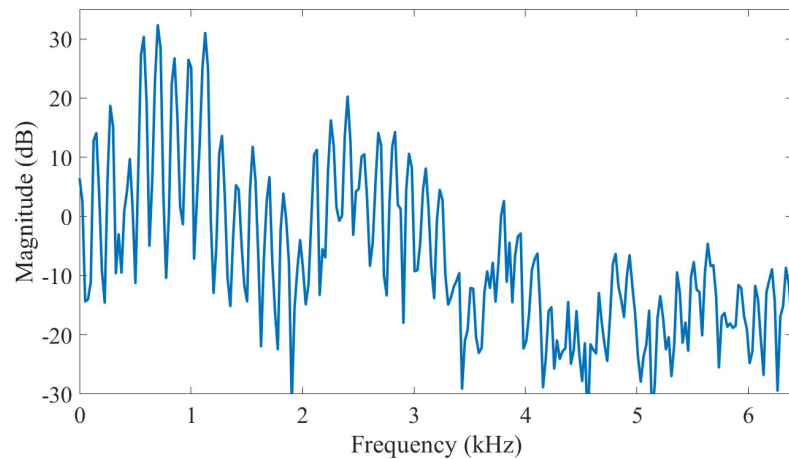# Visualising the cepstrum

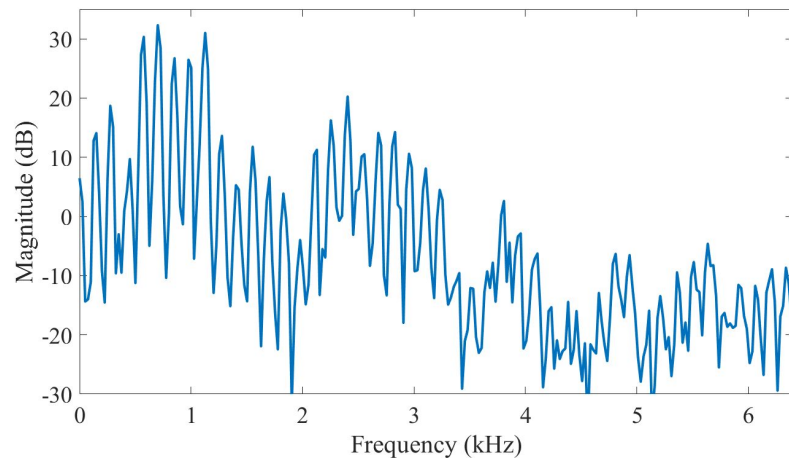Power spectrum



log

Log power spectrum

# Visualising the cepstrum

Log power spectrum

# Visualising the cepstrum

Log power spectrum

Cepstrum

IDFT

# Visualising the cepstrum

Log power spectrum
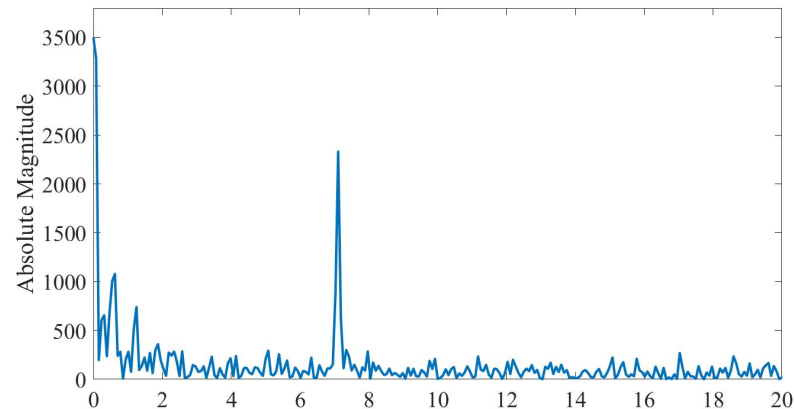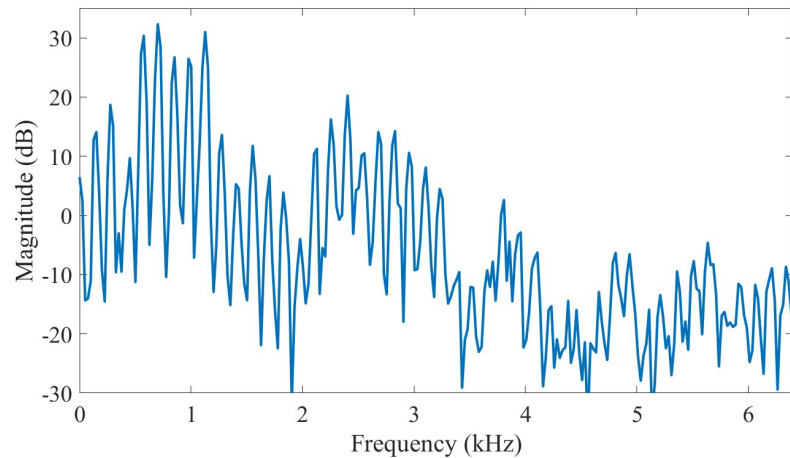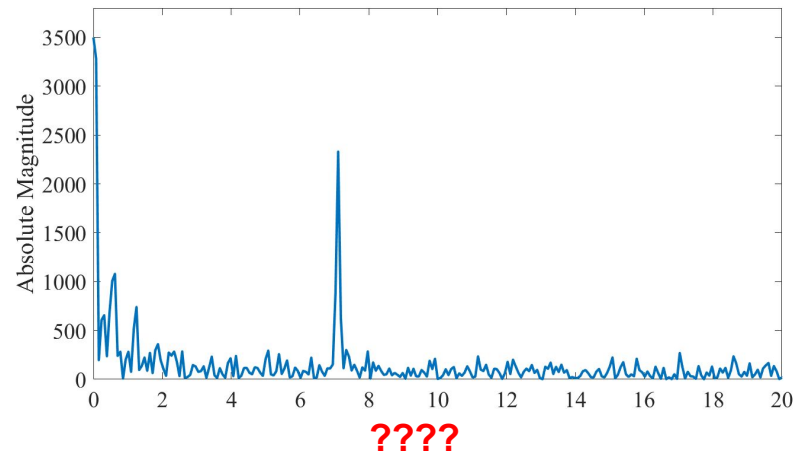
IDFT

Cepstrum



????

# Visualising the cepstrum

Log power spectrum

IDFT

Cepstrum
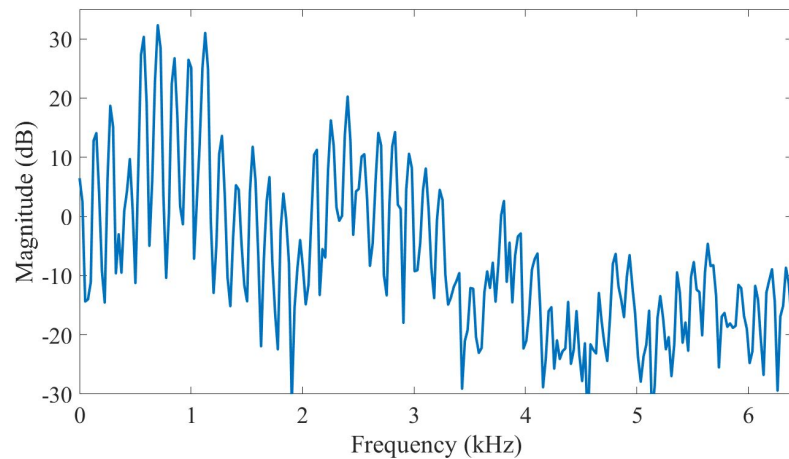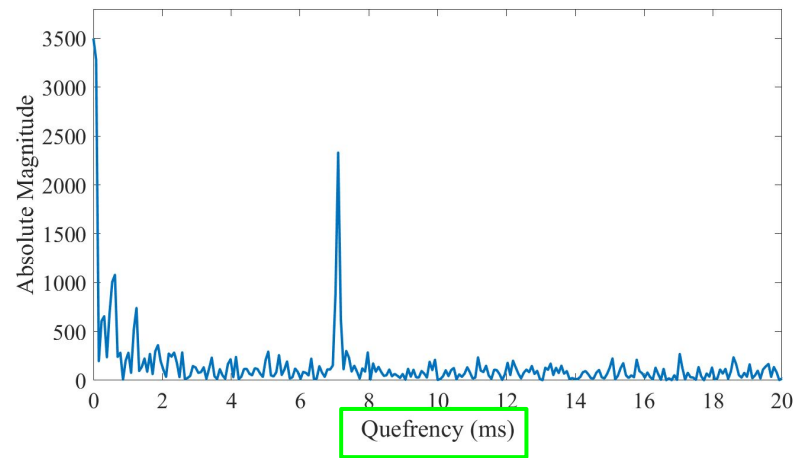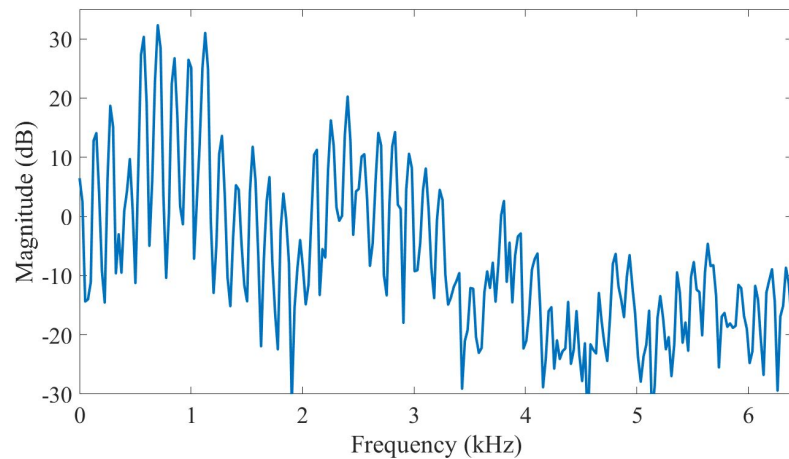
# Visualising the cepstrum



Log power spectrum
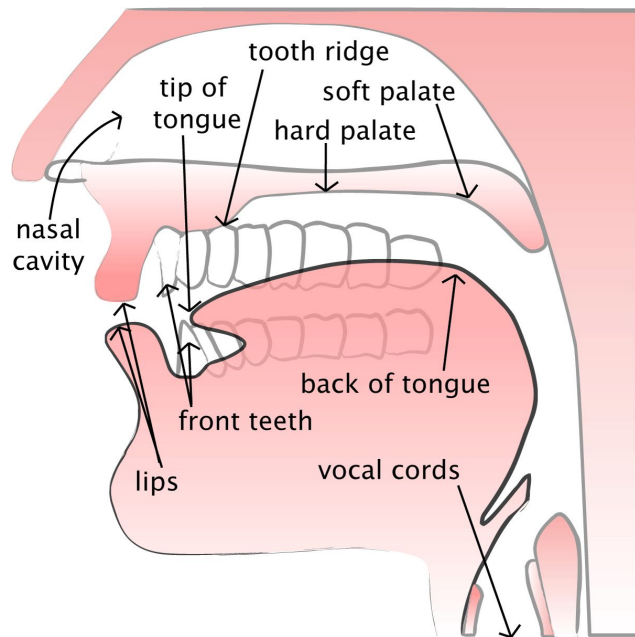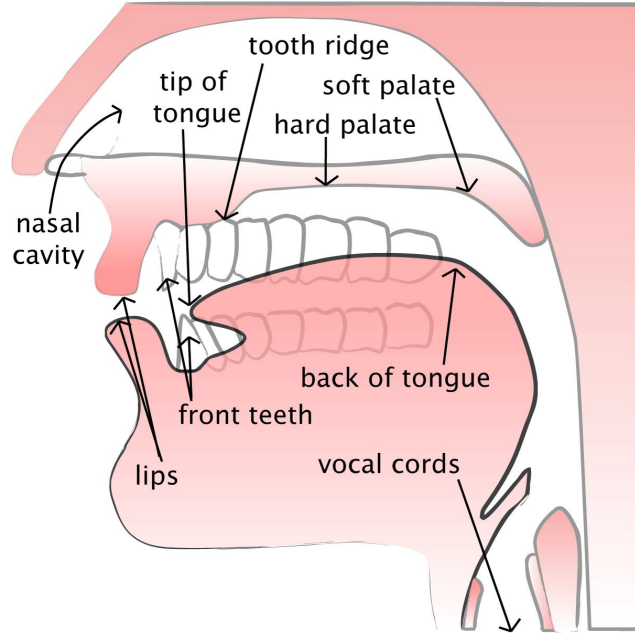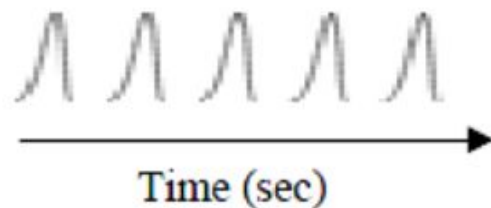
IDFT

Cepstrum

1st rhamonic

# The vocal tract

# The vocal tract



Vocal tract acts as a filter

# Speech generation



**Glottal pulses**

Time (sec)

**Vocal tract**

**Speech signal**

Time (sec)
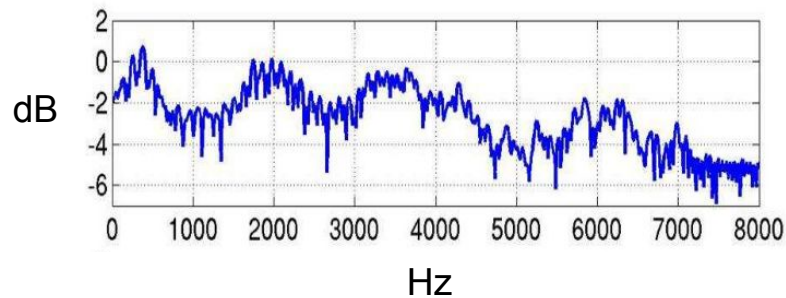
# Understanding the cepstrum
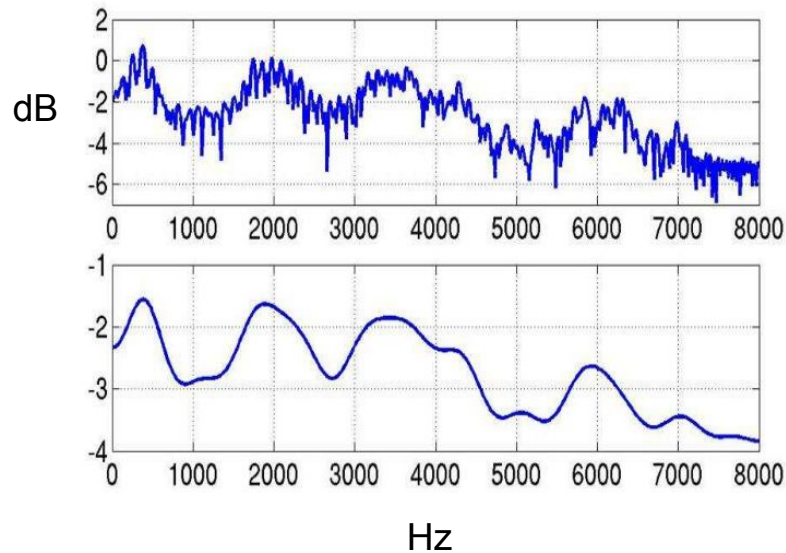
Log-spectrum



Speech

# Understanding the cepstrum

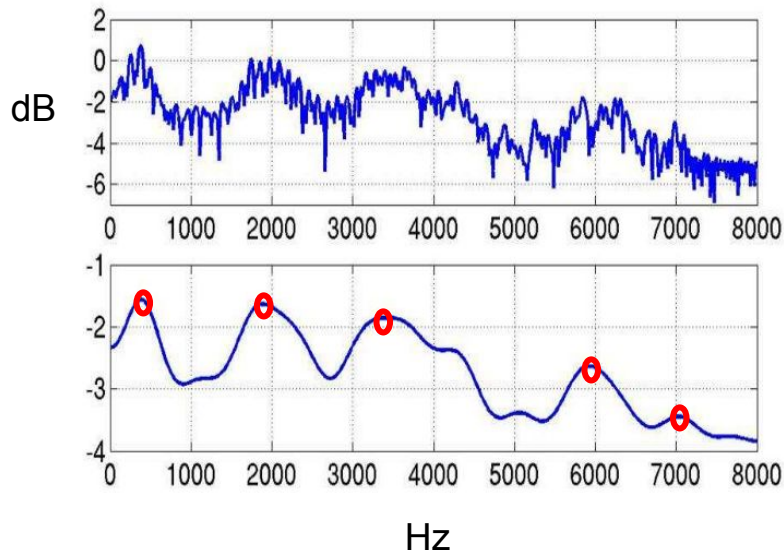Log-spectrum          dB                                    Speech



Spectral envelope

Hz
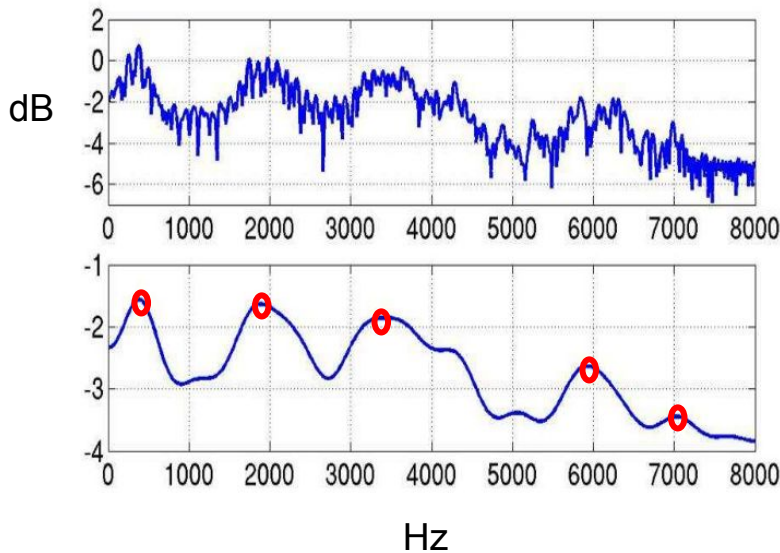
# Understanding the cepstrum

Log-spectrum

Spectral envelope

Speech

# Understanding the cepstrum
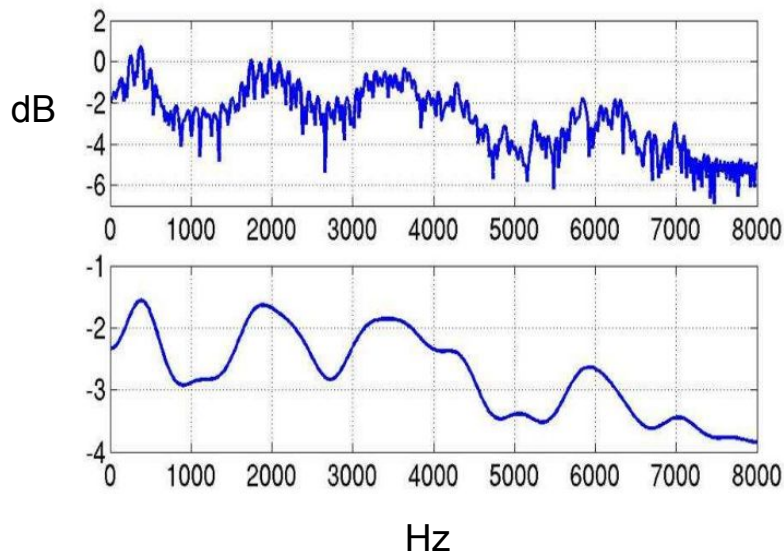
Log-spectrum

Spectral envelope

Speech



Formants = Carry identity of sound
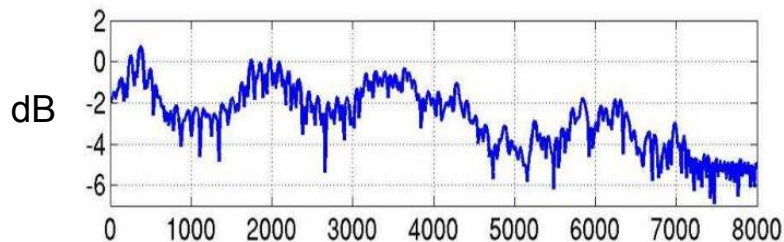
# Understanding the cepstrum

Log-spectrum

Spectral envelope



Speech

Vocal tract frequency response

# Understanding the cepstrum



Log-spectrum    dB

Spectral envelope

Hz

Speech

Vocal tract frequency response

# Understanding the cepstrum

Log-spectrum    dB

Spectral envelope

Spectral detail

Speech
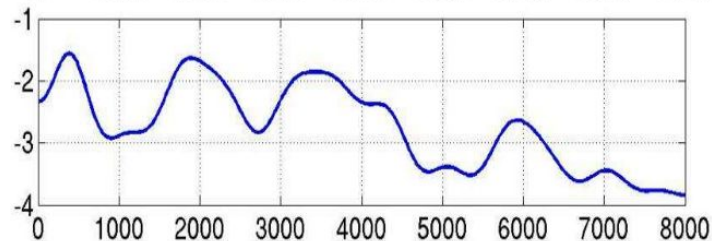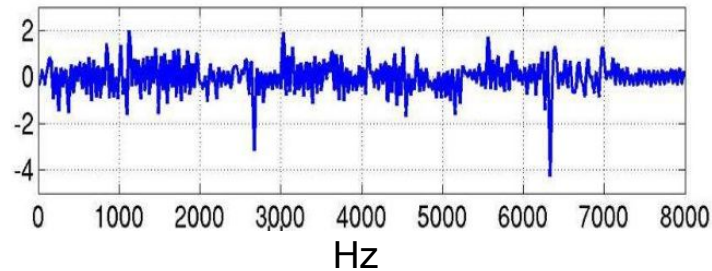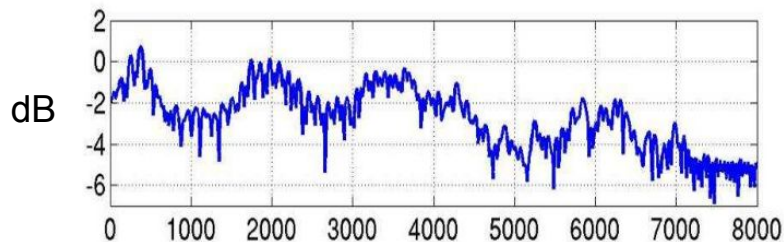
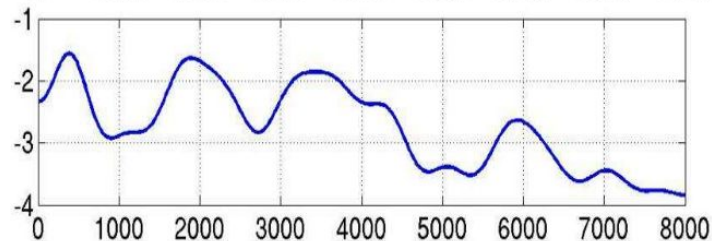Vocal tract frequency response

Hz

# Understanding the cepstrum

Log-spectrum  dB



Speech

Spectral envelope

Vocal tract frequency response

Spectral detail

Glottal pulse

Hz

# Speech

# =

# Convolution of vocal tract frequency response with glottal pulse

# Formalising speech

$$x(t) = e(t) \cdot h(t)$$

# Formalising speech

$$x(t) = e(t) \cdot h(t)$$

$$X(t) = E(t) \cdot H(t)$$

# Formalising speech

$$X(t) = E(t) \cdot H(t)$$

# Formalising speech

$$X(t) = E(t) \cdot H(t)$$

$$\Downarrow$$

$$log(X(t)) = log(E(t) \cdot H(t))$$

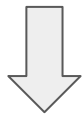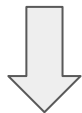# Formalising speech

$$X(t) = E(t) \cdot H(t)$$

$$\Downarrow$$

$$log(X(t)) = log(E(t) \cdot H(t))$$

$$\Downarrow$$
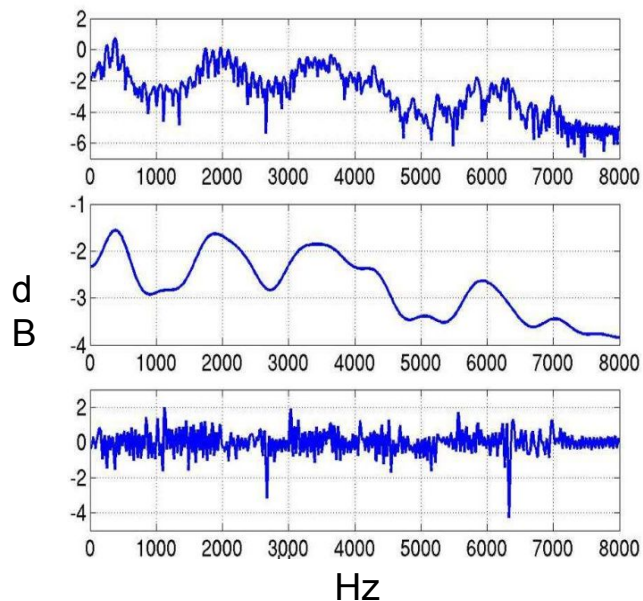
$$log(X(t)) = log(E(t)) + log(H(t))$$

# Formalising speech

$$log(X(t)) = log(E(t)) + log(H(t))$$

# Formalising speech

$$log(X(t)) = log(E(t)) + log(H(t))$$

Speech

Vocal tract frequency response

Glottal pulse

dB

Hz

# The goal: Separating components

# The goal: Separating components

# Separating components

$$log(X(t)) = log(E(t)) + log(H(t))$$

# Separating components

$$log(X(t)) = log(E(t)) + log(H(t))$$

IDFT

quefrency

d
B

Hz

# Separating components



$$log(X(t)) = log(E(t)) + log(H(t))$$

4 Hz

IDFT

quefrency

d
B

Hz

# Separating components

$$log(X(t)) = log(E(t)) + log(H(t))$$



100 Hz

quefrency

IDFT

d
B

Hz

# Separating components



$$log(X(t)) = log(E(t)) + log(H(t))$$

quefrency

IDFT

dB

Hz

$$X(t) = \boxed{E(t)} + \boxed{H(t)}$$

# Separating components



$$log(X(t)) = log(E(t)) + log(H(t))$$

quefrency

IDFT
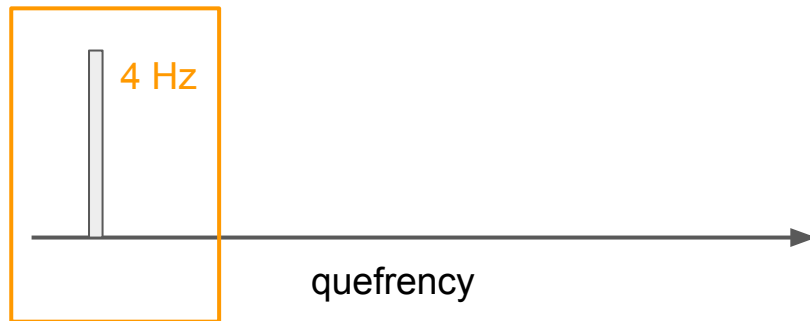
dB

Hz

$$X(t) = E(t) + H(t)$$

# Computing Mel-Frequency Cepstral Coefficients

Waveform

# Computing Mel-Frequency Cepstral Coefficients

```
┌─────────────────┐
│    Waveform     │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│       DFT       │
└─────────────────┘
```

# Computing Mel-Frequency Cepstral Coefficients

Waveform

↓

DFT

↓

Log-Amplitude
Spectrum

# Computing Mel-Frequency Cepstral Coefficients

Waveform

DFT

Log-Amplitude Spectrum

Mel-Scaling

# Computing Mel-Frequency Cepstral Coefficients

Waveform

↓

DFT

↓

Log-Amplitude
Spectrum

↓

Mel-Scaling

↓

Discrete
Cosine
Transform

↓

MFCCs

# Why Discrete Cosine Transform?

# Why Discrete Cosine Transform?

- Simplified version of Fourier Transform

# Why Discrete Cosine Transform?

- Simplified version of Fourier Transform

- Get real-valued coefficient

# Why Discrete Cosine Transform?

- Simplified version of Fourier Transform

- Get real-valued coefficient

# Why Discrete Cosine Transform?

- Simplified version of Fourier Transform

- Get real-valued coefficient

- Decorrelate energy in different mel bands

# Why Discrete Cosine Transform?

- Simplified version of Fourier Transform

- Get real-valued coefficient
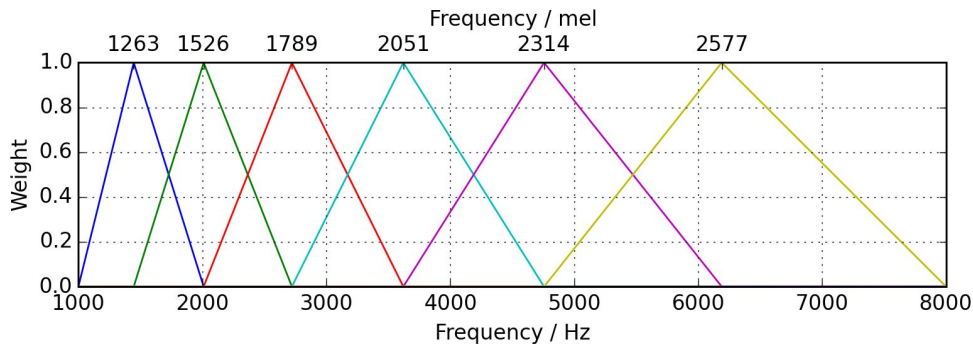
- Decorrelate energy in different mel bands

- Reduce # dimensions to represent spectrum

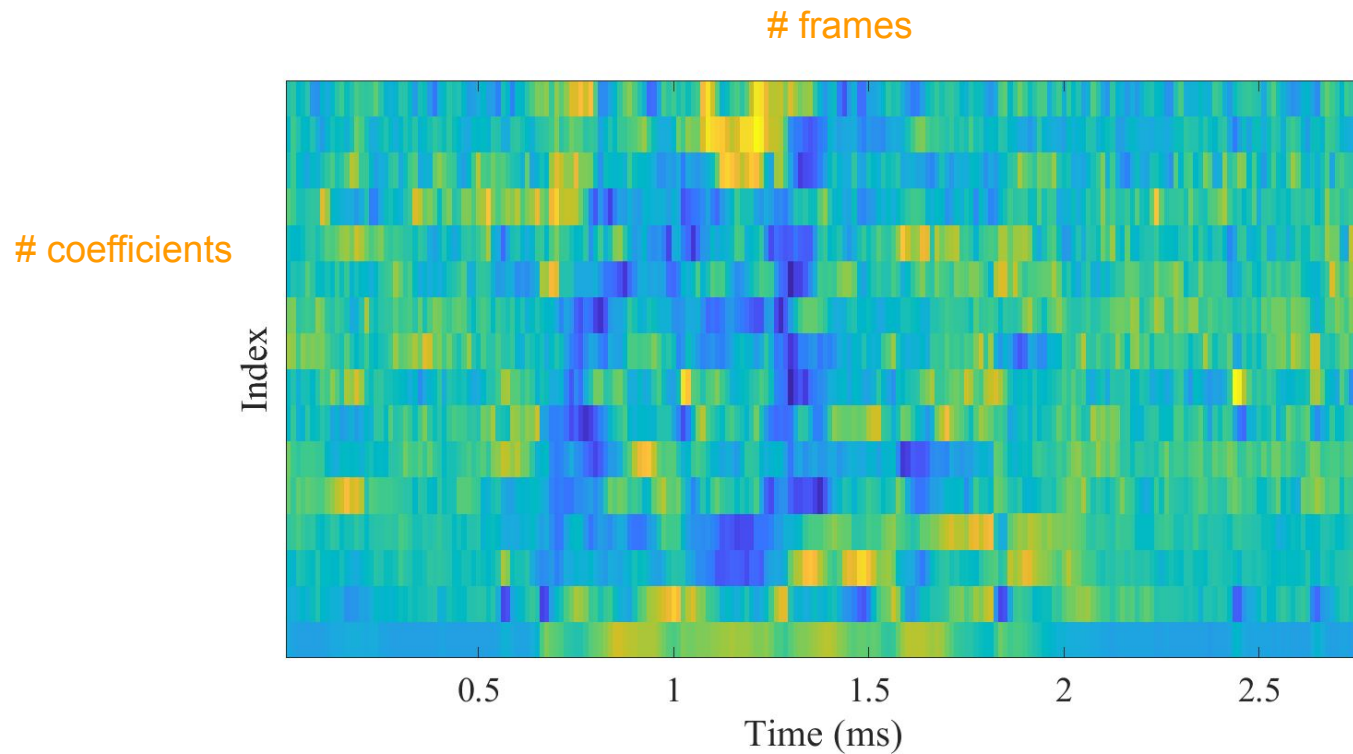# How many coefficients?

- Traditionally: first 12 - 13 coefficients

- First coefficients keep most information (e.g., formants, spectral envelope)

- Use Δ and ΔΔ MFCCs

- Total 39 coefficients per frame

# Visualising MFCCs

# Visualising MFCCs

# MFCCs advantages

- Describe the "large" structures of the spectrum

- Ignore fine spectral structures

- Work well in speech and music processing

# MFCCs disadvantages

- Not robust to noise

- Extensive knowledge engineering

- Not efficient for synthesis

# MFCCs applications

- Speech processing

  - Speech recognition

  - Speaker recognition

  - ...

- Music processing

  - Music genre classification

  - Mood classification

  - Automatic tagging

  - ...

# What's up next?

- Extract MFCCs with Python and Librosa

- Visualise MFCCs