# Mel-Frequency Cepstral Coefficients

Cepstrum

⇩

Spectrum

| 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

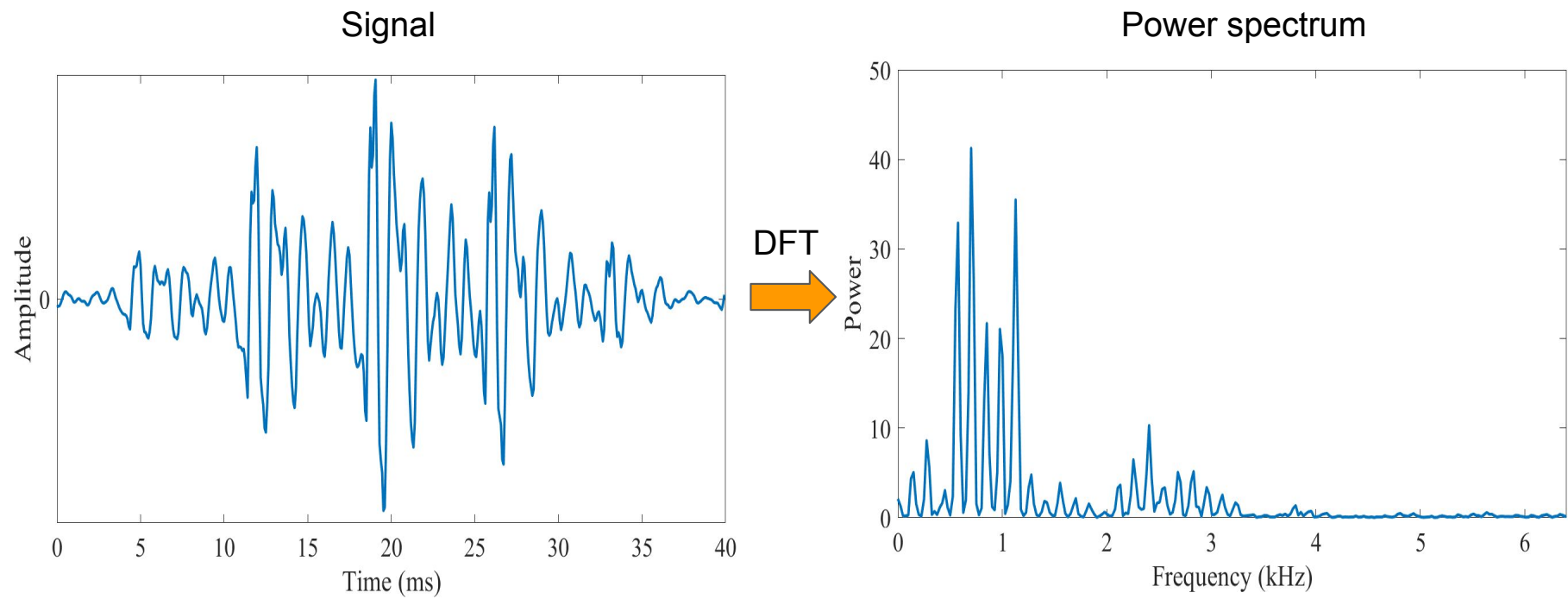$$C(\underbrace{x(t)}_{\text{Time-domain signal}}) = \underbrace{F^{-1}[\underbrace{log(\underbrace{F[\underbrace{x(t)}]}_{\text{Spectrum}})}_{\text{Log spectrum}}]}_{\text{Cepstrum}}$$

Time-domain signal

Spectrum

Log spectrum

Cepstrum

# Visualising the cepstrum



Signal

Power spectrum

DFT

# Visualising the cepstrum
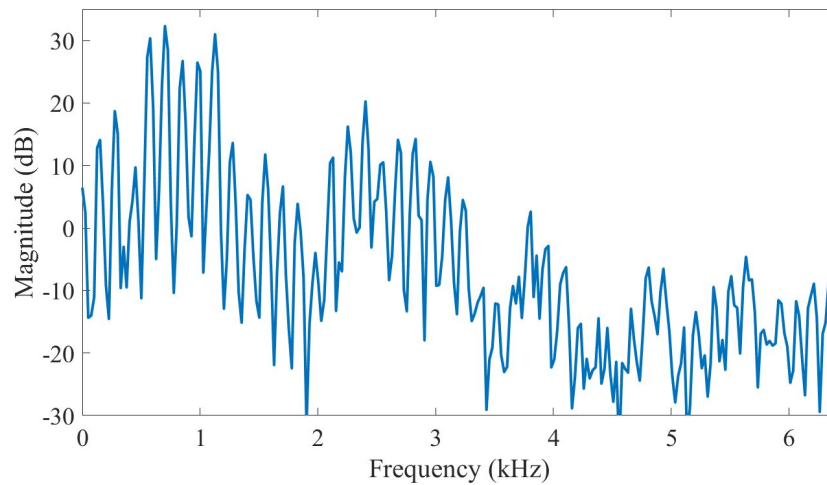


Power spectrum

log

Log power spectrum

# Visualising the cepstrum

Log power spectrum

IDFT

Cepstrum

# Visualising the cepstrum



Log power spectrum

IDFT

Cepstrum

????

# Visualising the cepstrum

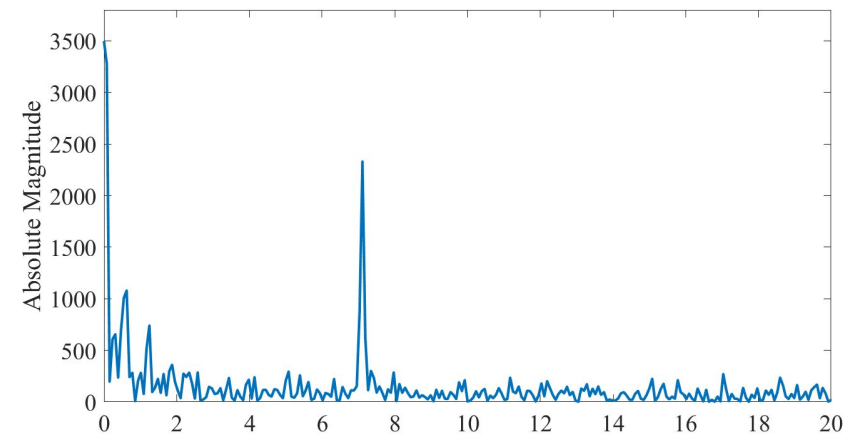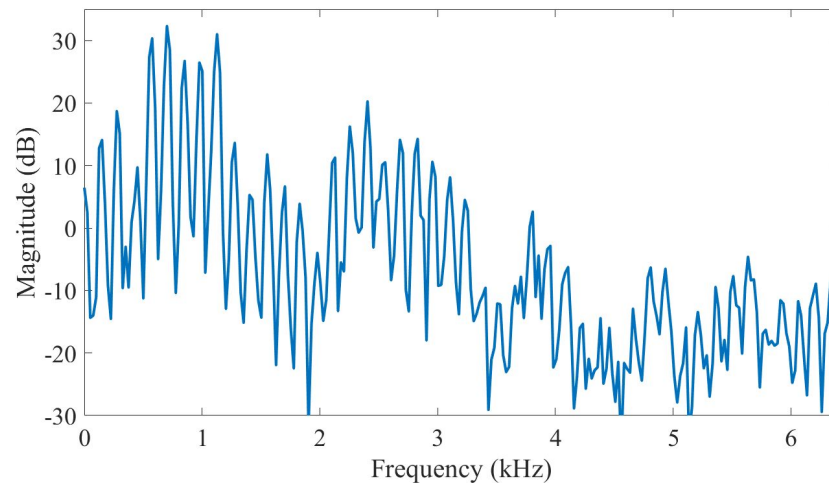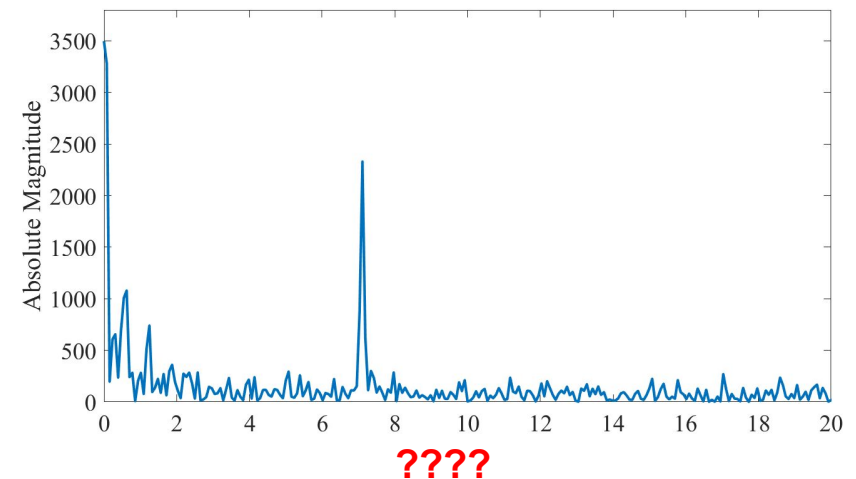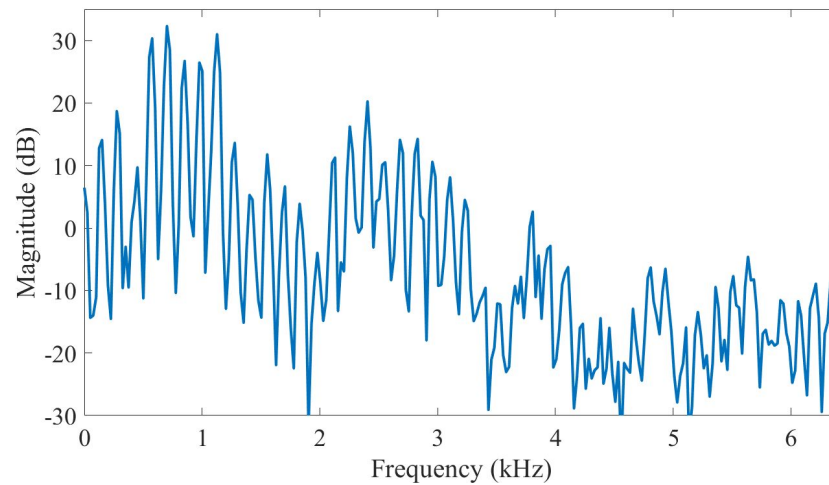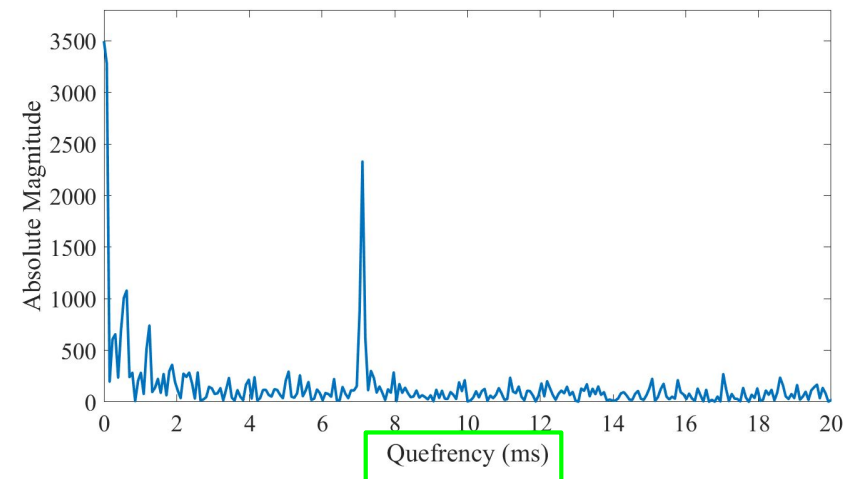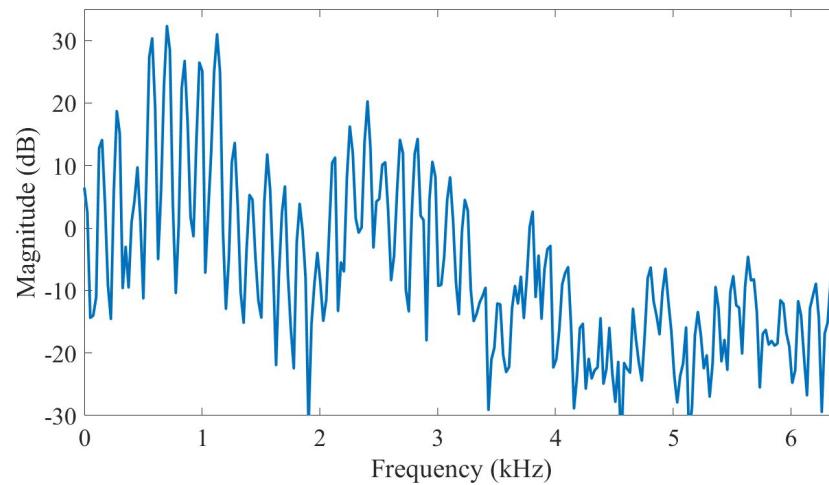Log power spectrum
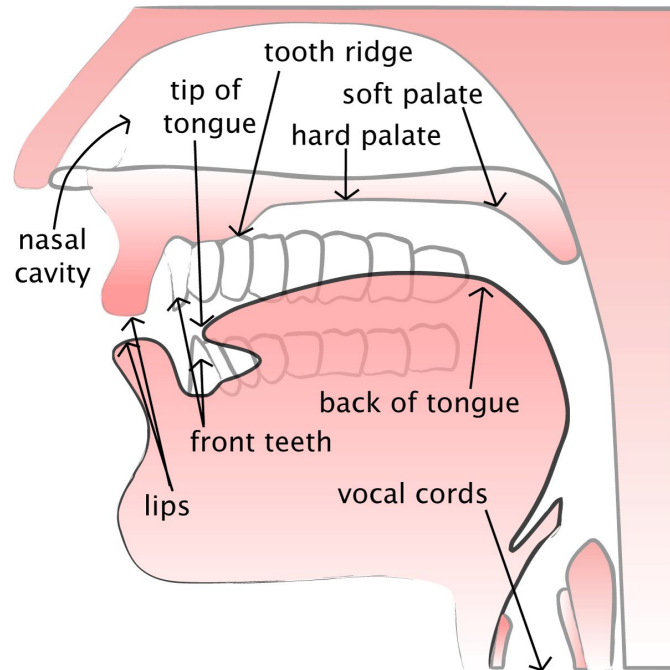
IDFT

Cepstrum
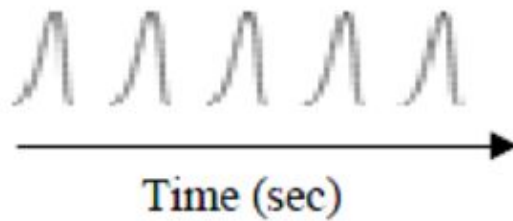
# Visualising the cepstrum

Log power spectrum

IDFT

Cepstrum

1st rhamonic

# The vocal tract



Vocal tract acts as a filter

# Speech generation



Glottal pulses     Vocal tract     Speech signal

Time (sec)                 Time (sec)

# Understanding the cepstrum

Log-spectrum



Speech

# Understanding the cepstrum

Log-spectrum          dB

Spectral envelope

Hz
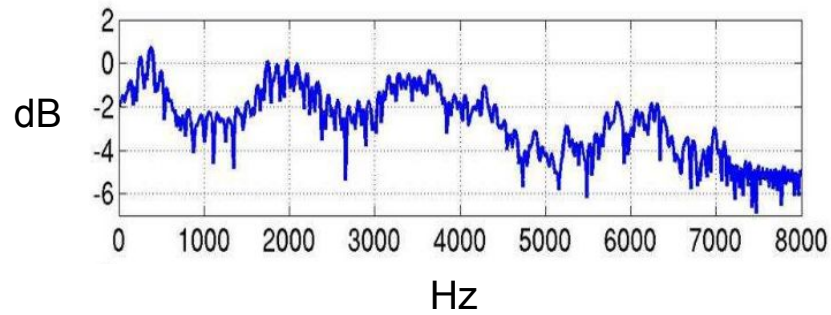
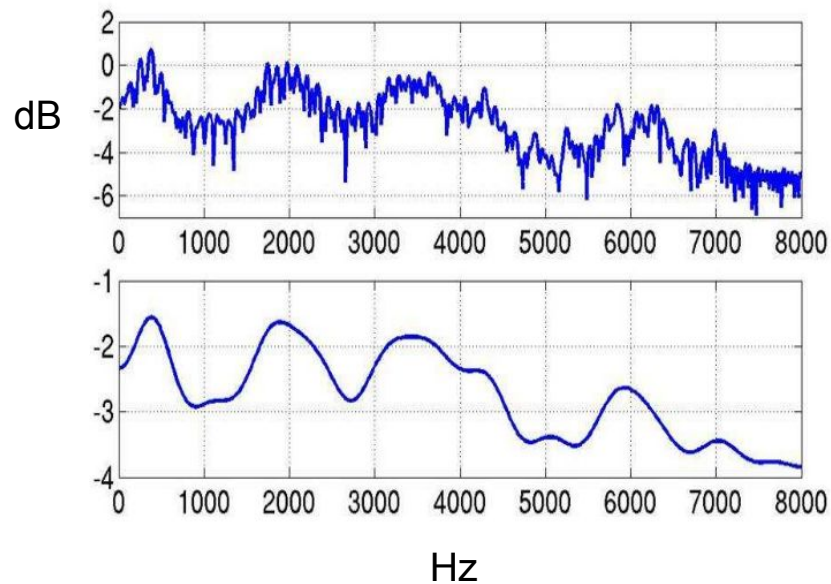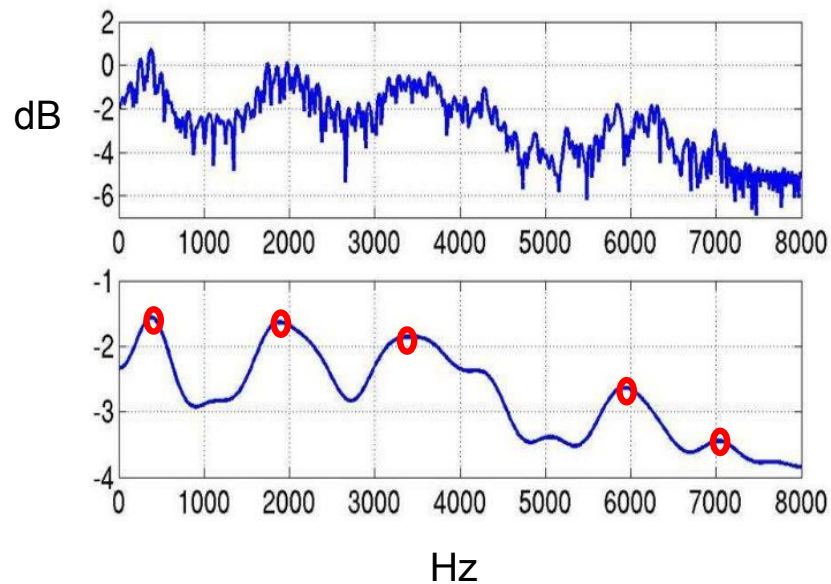Speech

# Understanding the cepstrum



Log-spectrum    dB    Speech
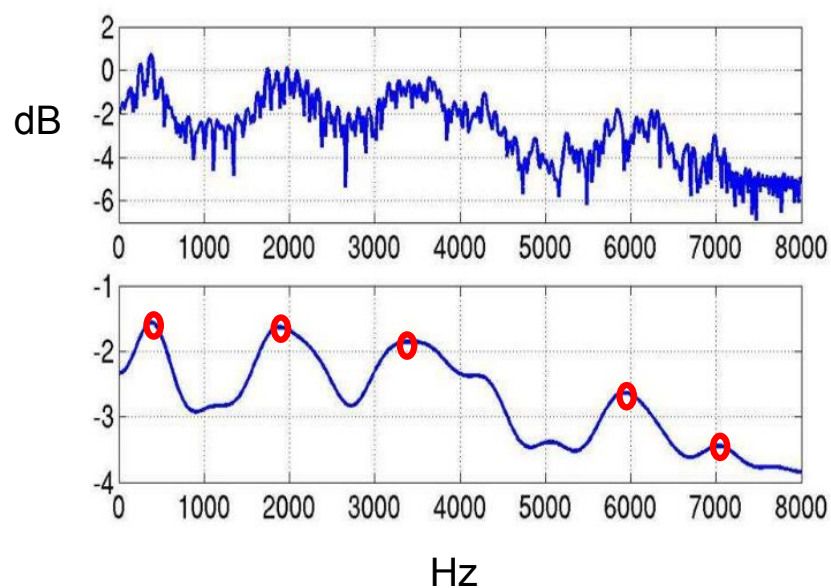
Spectral envelope

Hz

# Understanding the cepstrum



Log-spectrum     dB                                           Speech

Spectral envelope

Hz

**Formants = Carry identity of sound**

# Understanding the cepstrum



Log-spectrum    dB                                               Speech

Spectral envelope                                      Vocal tract frequency response

Hz

# Understanding the cepstrum



Log-spectrum    dB                            Speech

Spectral envelope                        Vocal tract frequency response

Hz

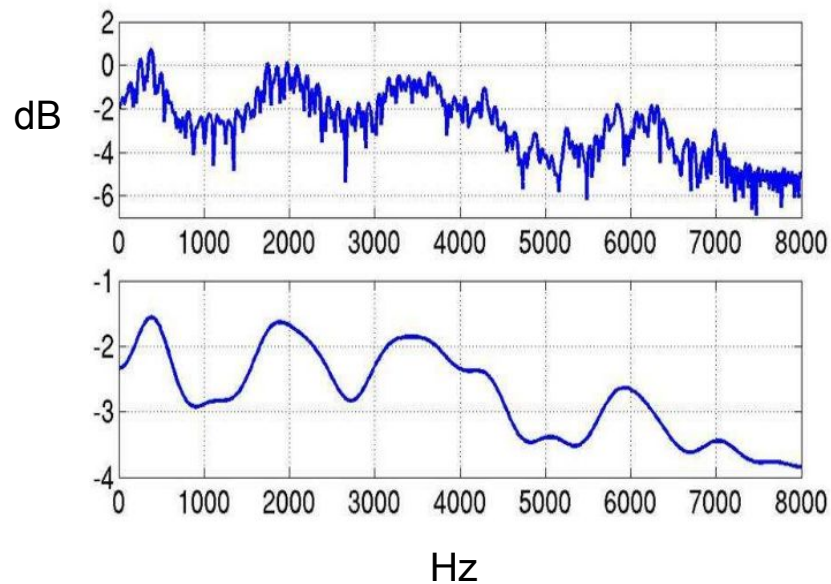# Understanding the cepstrum

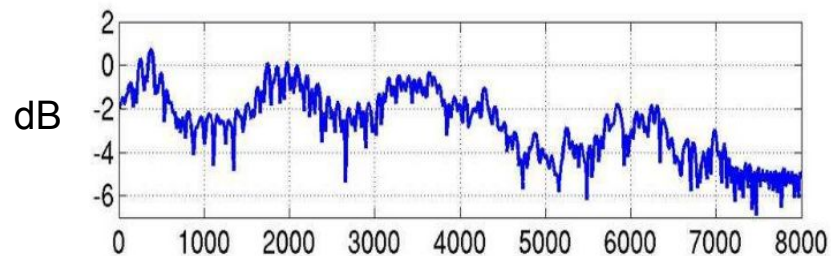Log-spectrum    dB                                              Speech

Spectral envelope                                              Vocal tract frequency
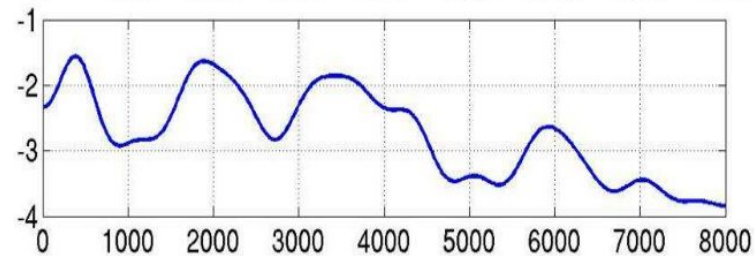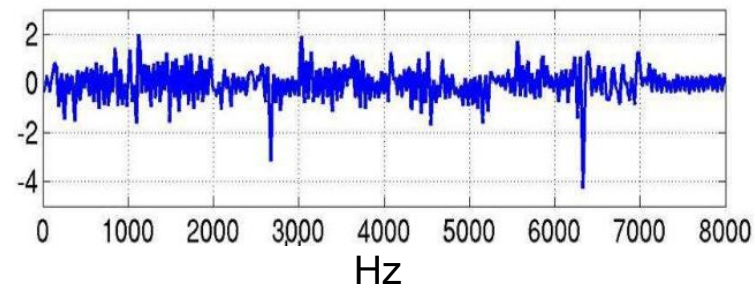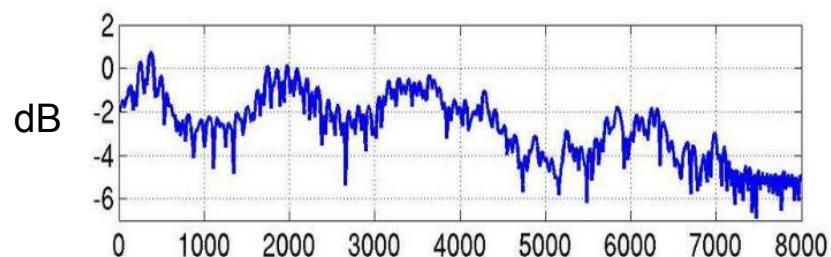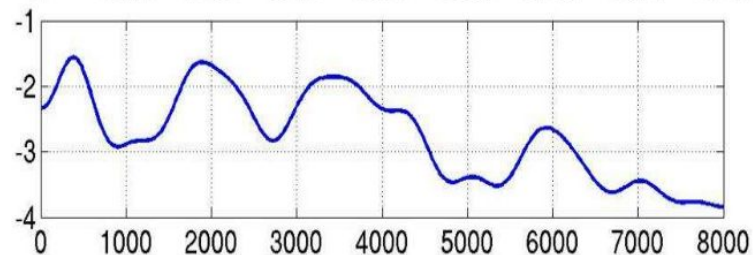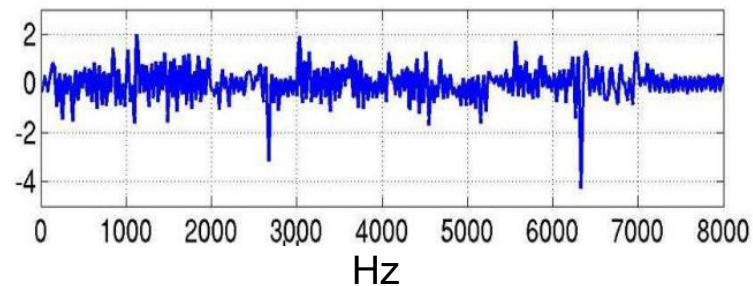                                                               response
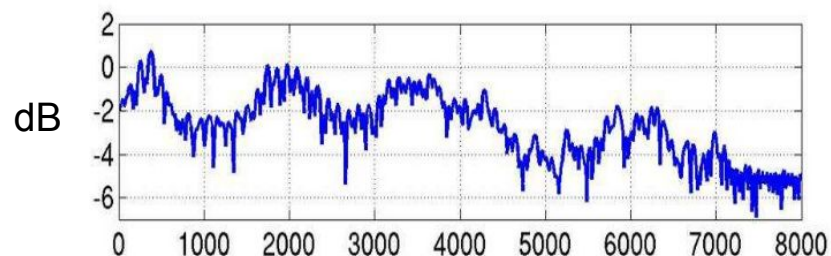
Spectral detail

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

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

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

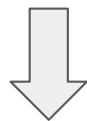## Formalising speech

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

$$\Downarrow$$

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

$$\Downarrow$$

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

# 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

d
B

Hz

# The goal: Separating components

# The goal: Separating components

# Separating components

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



dB

Hz

# Separating components

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

quefrency

IDFT

d
B

Hz

# Separating components



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

4 Hz

quefrency

IDFT

d
B

Hz

# Separating components



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

quefrency

100 Hz

IDFT

d
B

Hz

# Separating components



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

IDFT

quefrency

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

d
B

Hz

# Separating components



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

quefrency

IDFT

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

dB

Hz

# Computing Mel-Frequency Cepstral Coefficients

```
┌─────────────────┐
│    Waveform     │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│      DFT        │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│  Log-Amplitude  │
│    Spectrum     │
└─────────────────┘
```

# Computing Mel-Frequency Cepstral Coefficients

```
Waveform
   |
   v
  DFT
   |
   v
Log-Amplitude
  Spectrum
   |
   v
Mel-Scaling
```

# Computing Mel-Frequency Cepstral Coefficients

```
┌─────────────────┐
│    Waveform     │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│       DFT       │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│  Log-Amplitude  │
│    Spectrum     │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│   Mel-Scaling   │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│    Discrete     │
│     Cosine      │
│    Transform    │
└─────────────────┘
         │
         ▼
       MFCCs
```

# 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
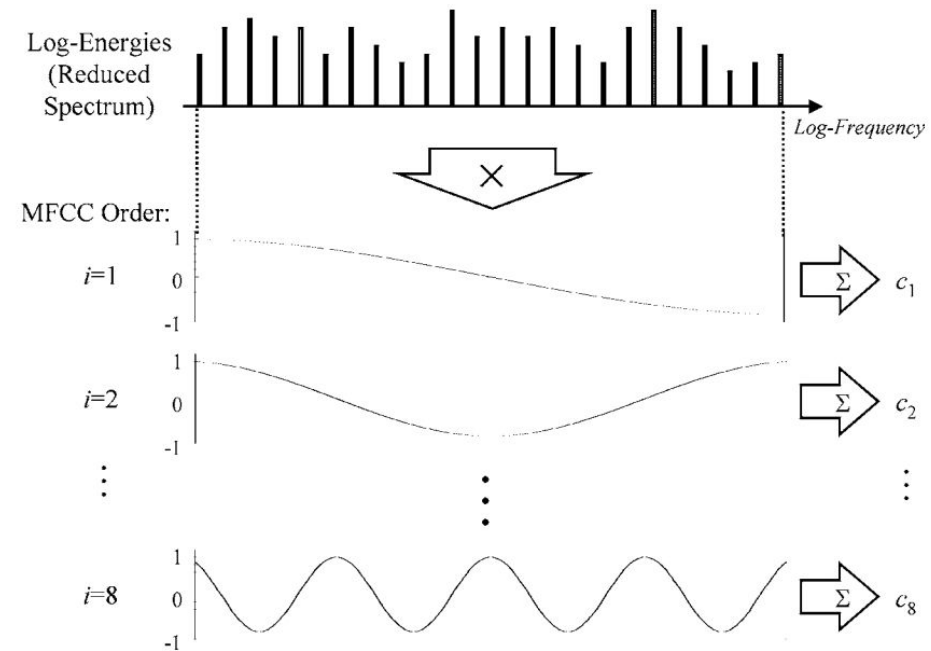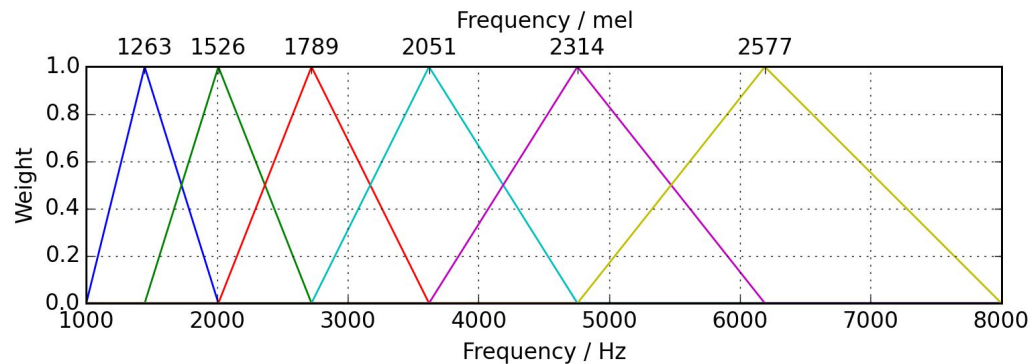
# Why Discrete Cosine Transform?
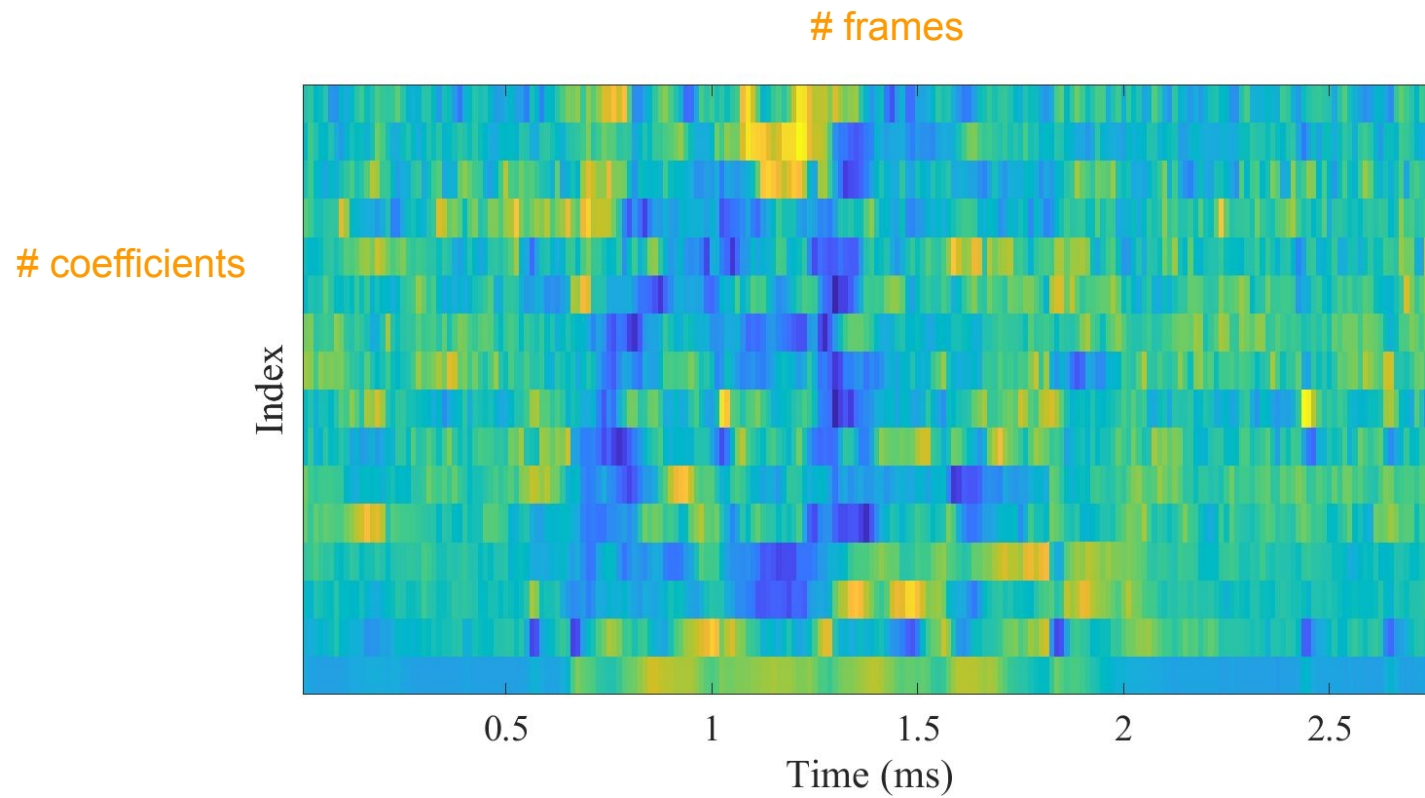
- Simplified version of Fourier Transform

- 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

# 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