

Speech Signal Analysis

Hiroshi Shimodaira and Steve Renals

Automatic Speech Recognition— ASR Lectures 2&3
17/24 January 2013

Overview

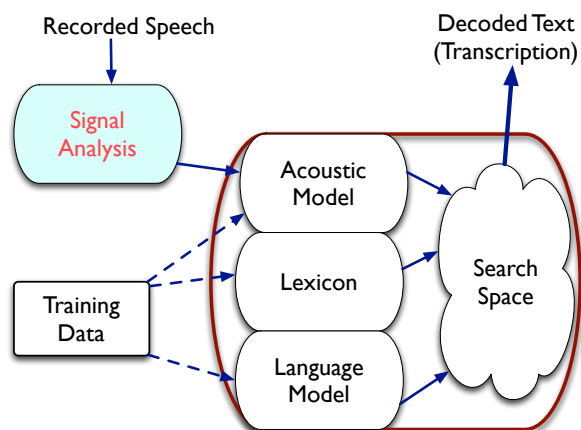
Speech Signal Analysis for ASR

- Features for ASR
- Spectral analysis
- Cepstral analysis
- Standard features for ASR: MFCCs and PLP analysis
- Dynamic features

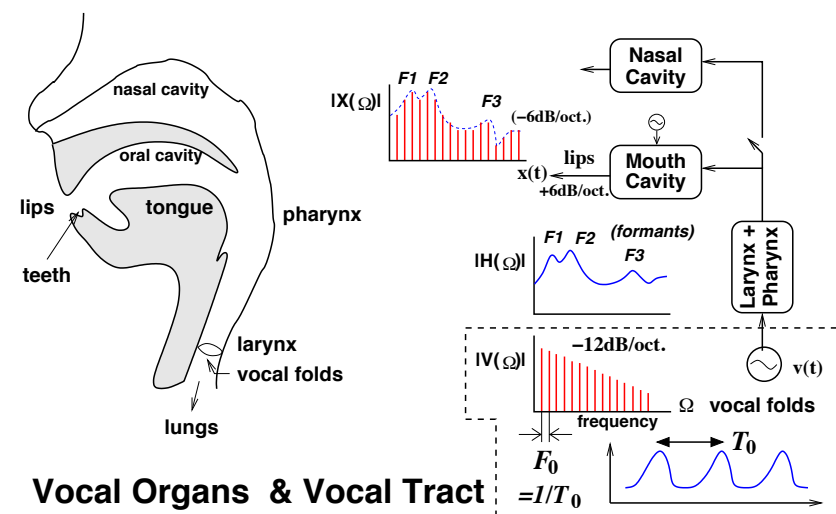
Reading:

- Jurafsky & Martin, sec 9.3
- P Taylor, *Text-to-Speech Synthesis*, chapter 12, signal processing background chapter 10

Speech signal analysis for ASR

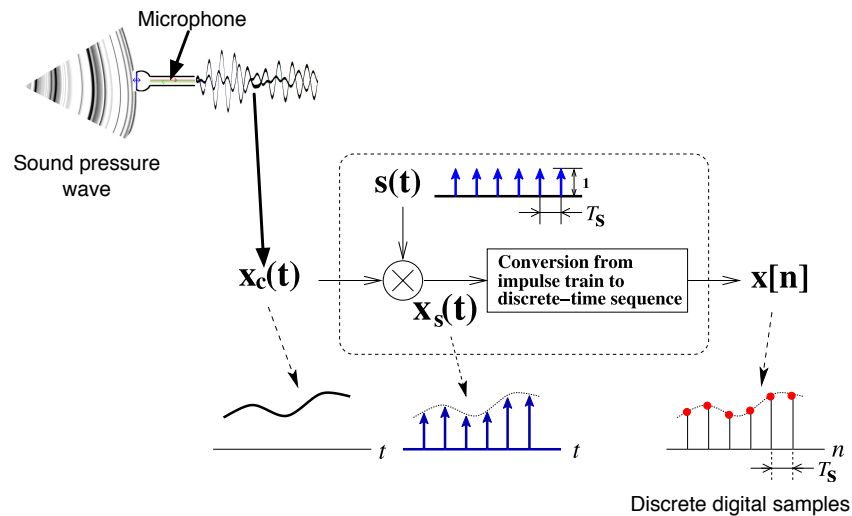


Speech production model

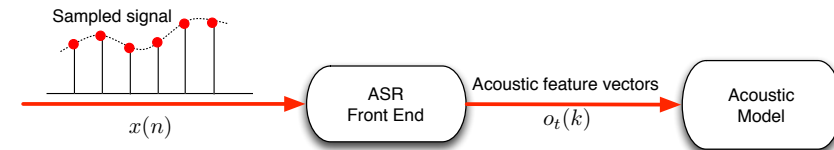


Vocal Organs & Vocal Tract

Sampling



Acoustic Features for ASR



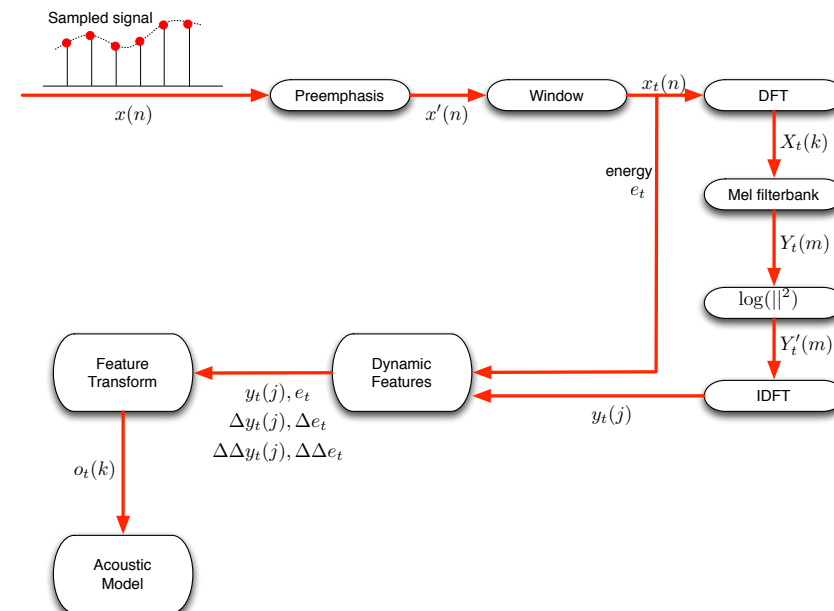
Speech signal analysis to produce a sequence of acoustic feature vectors

Acoustic Features for ASR

Desirable characteristics of acoustic features used for ASR:

- Features should contain sufficient information to distinguish between phones
 - good time resolution (10ms)
 - good frequency resolution (~ 20 channels)
- Be separated from F_0 and its harmonics
- Be robust against speaker variation
- Be robust against noise or channel distortions
- Have good "pattern recognition characteristics"
 - low feature dimension
 - features are independent of each other

MFCC-based front end for ASR

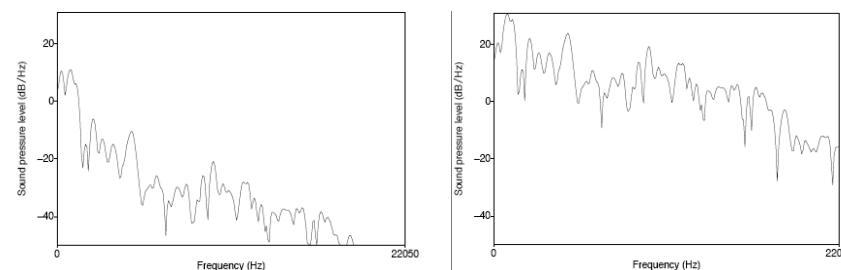


Pre-emphasis and spectral tilt

- Pre-emphasis increases the magnitude of higher frequencies in the speech signal compared with lower frequencies
- *Spectral Tilt*
 - The speech signal has more energy at low frequencies (for voiced speech)
 - This is due to the glottal source
- Pre-emphasis (first-order) filter boosts higher frequencies:

$$x[n'] = x[n] + \alpha x[n - 1] \quad 0.95 < \alpha < 0.99$$

Pre-emphasis: example



Vowel /aa/ - time slice of the spectrum

(Jurafsky & Martin, fig. 9.9)

Windowing

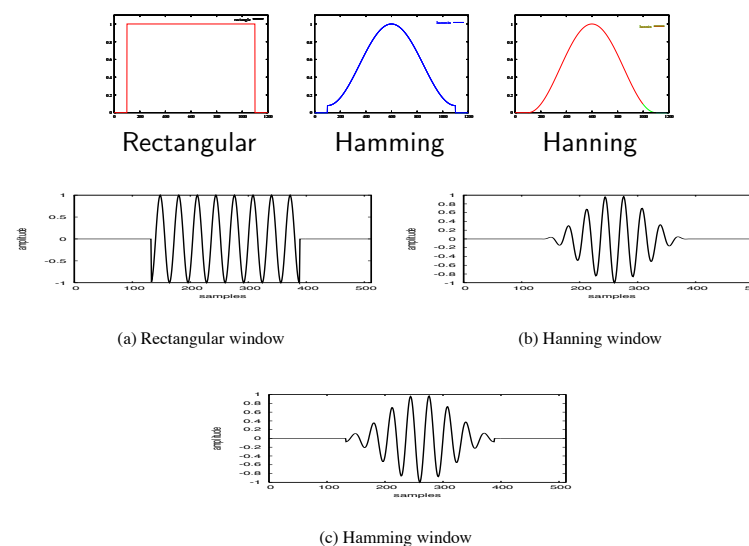
- The speech signal is constantly changing (non-stationary)
- Signal processing algorithms usually assume that they the signal is stationary
- Piecewise stationarity: model speech signal as a sequence of **frames** (each assumed to be stationary)
- **Windowing**: multiply the full waveform $s(n)$ by a window $w(n)$ (in time domain):

$$x[n] = w[n]s[n]$$

- Simply cutting out a short segment (frame) from $s(n)$ is a rectangular window — causes discontinuities at the edges of the segment
- Instead, a tapered window is usually used
e.g. *Hamming* ($\alpha = 0.46164$) or *Hanning* ($\alpha = 0.5$) window

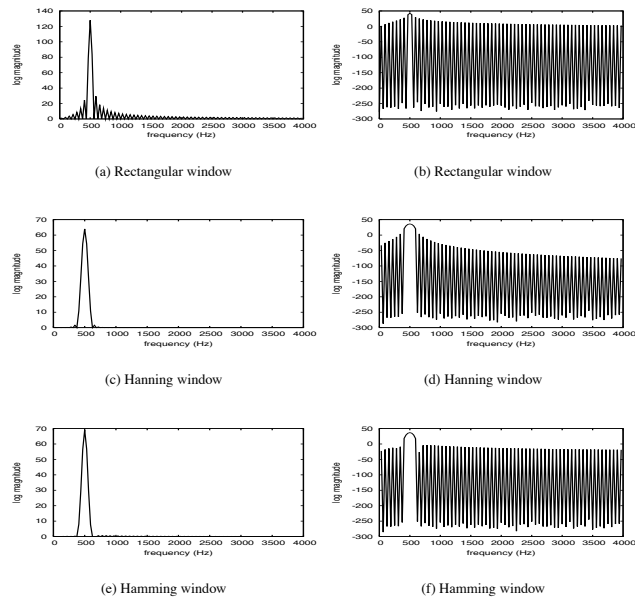
$$w[\ell] = (1 - \alpha) - \alpha \cos\left(\frac{2\pi\ell}{N-1}\right) \quad N : \text{window width}$$

Effect of windowing — time domain



(Taylor, fig 12.1)

Effect of windowing — frequency domain



(Taylor, fig 12.2)

Discrete Fourier Transform (DFT)

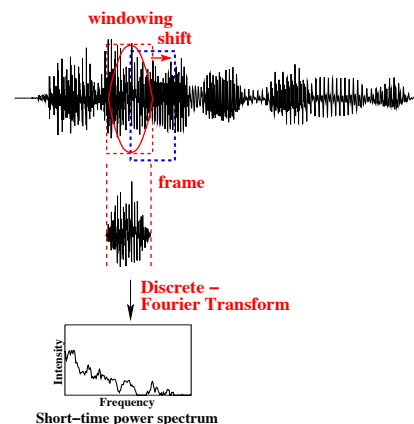
- Purpose: extracts spectral information from a windowed signal (i.e. how much energy at each frequency band)
- Input: windowed signal $x[n] \dots x[m]$ (time domain)
- Output: a complex number $X[k]$ for each of N frequency bands representing magnitude and phase for the k th frequency component (frequency domain)
- Discrete Fourier Transform (DFT):

$$X[k] = \sum_{n=0}^{N-1} x[n] \exp\left(-j \frac{2\pi}{N} kn\right)$$

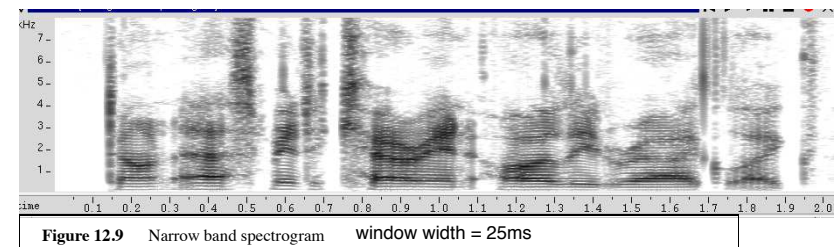
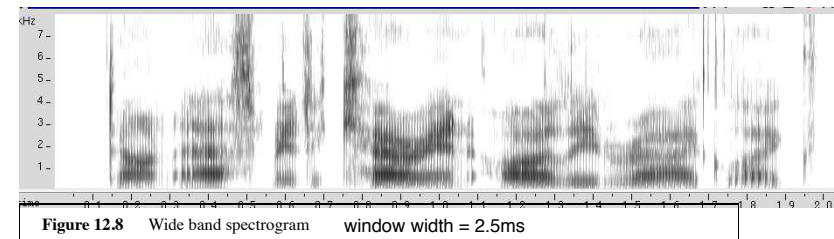
- Fast Fourier Transform (FFT) — efficient algorithm for computing DFT when N is a power of 2

Windowing and spectral analysis

- Window the signal $x(n)$ into frames $x_t(n)$ and apply Fourier Transform to each segment.
 - Short frame width: *wide-band*, high time resolution, low frequency resolution
 - Long frame width: *narrow-band*, low time resolution, high frequency resolution
- For ASR:
 - frame width $\sim 20ms$
 - frame shift $\sim 10ms$

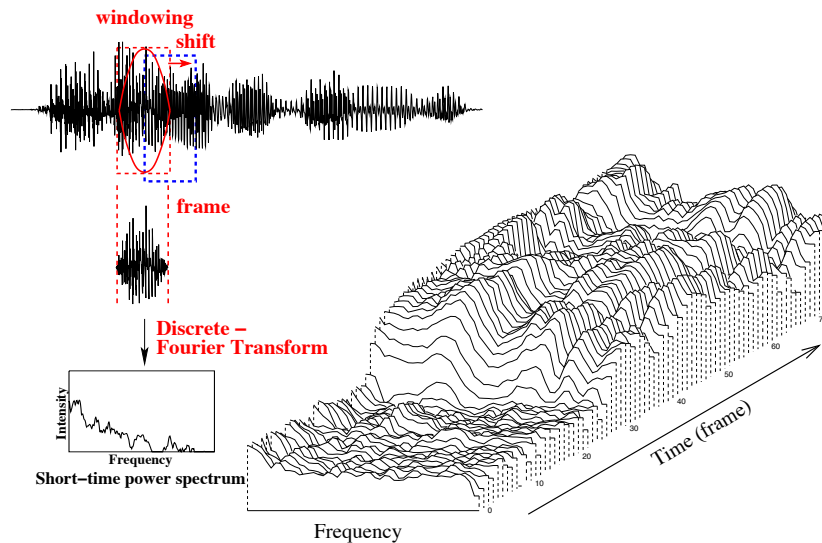


Wide-band and narrow-band spectrograms

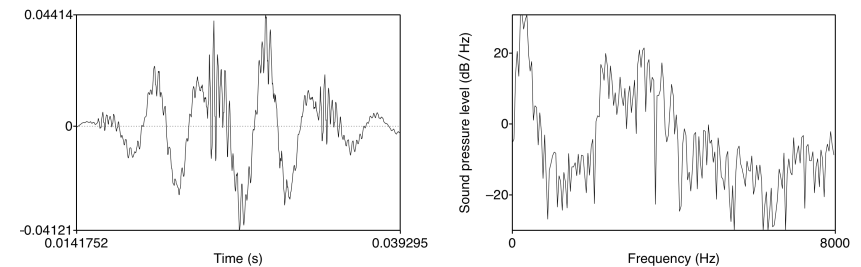


(Taylor, figs 12.8, 12.9)

Short-time spectral analysis



DFT Spectrum

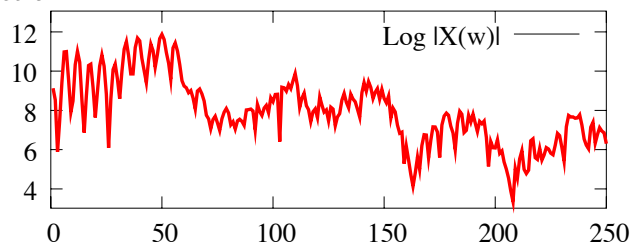


25ms Hamming window of vowel /iy/ and its spectrum computed by DFT

(Jurafsky and Martin, fig 9.12)

DFT Spectrum Features for ASR

- Equally-spaced frequency bands — but human hearing less sensitive at higher frequencies (above $\sim 1000\text{Hz}$)
- The estimated power spectrum contains harmonics of F_0 , which makes it difficult to estimate the envelope of the spectrum



- Frequency bins of STFT are highly correlated each other, i.e. power spectrum representation is highly redundant

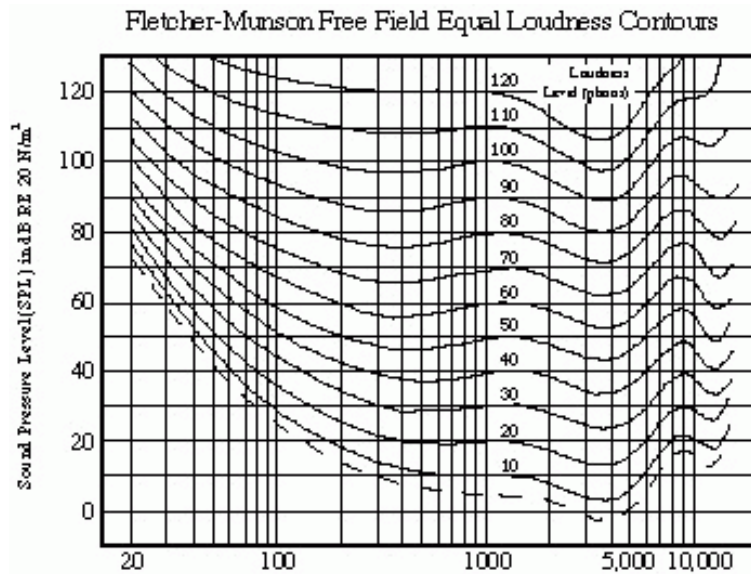
Human hearing

Physical quality	Perceptual quality
Intensity	Loudness
Fundamental frequency	Pitch
Spectral shape	Timbre
Onset/offset time	Timing
Phase difference in binaural hearing	Location

Technical terms

- equal-loudness contours
- masking
- auditory filters (critical-band filters)
- critical bandwidth

Equal loudness contour



ASR Lectures 2&3

Speech Signal Analysis

21

Nonlinear frequency scaling

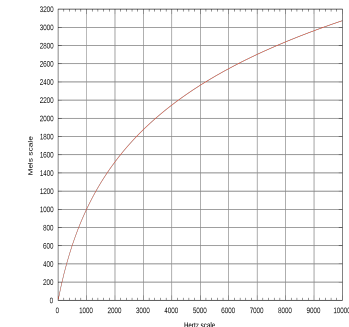
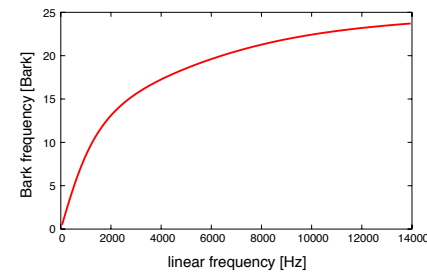
Human hearing is less sensitive to higher frequencies — thus human perception of frequency is nonlinear

Bark scale

Mel scale

$$b(f) = 13 \arctan(0.00076f) + 3.5 \arctan((f/7500)^2)$$

$$M(f) = 1127 \ln(1 + f/700)$$



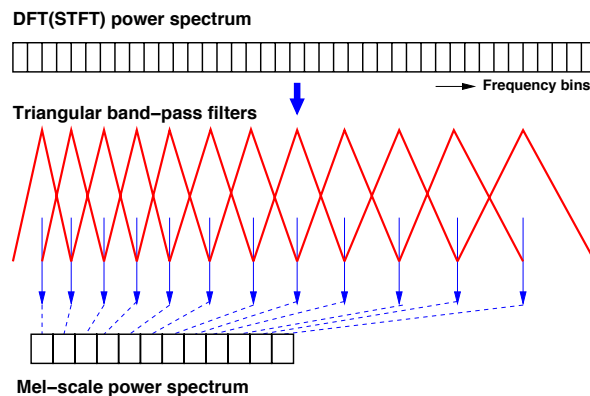
ASR Lectures 2&3

Speech Signal Analysis

22

Mel Filterbank

- Apply a mel-scale filter bank to DFT power spectrum to obtain mel-scale power spectrum
- Each filter collects energy from a number of frequency bands in the DFT
- Linearly spaced < 1000 Hz, logarithmically spaced > 1000 Hz



ASR Lectures 2&3

Speech Signal Analysis

23

Log Energy

- Compute the log magnitude squared of each Mel filter bank output
 - Taking the log compresses the dynamic range
 - Human sensitivity to signal energy is logarithmic — i.e. humans are less sensitive to small changes in energy at high energy than small changes at low energy
 - Log makes features less variable to acoustic coupling variations
 - Removes phase information — not important for speech recognition (not everyone agrees with this)

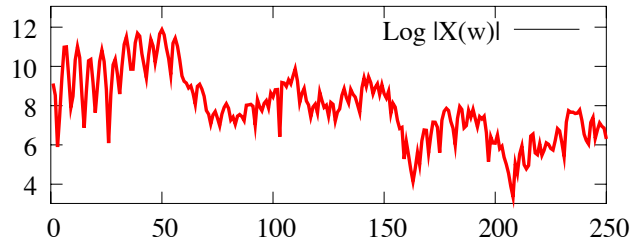
ASR Lectures 2&3

Speech Signal Analysis

24

DFT Spectrum Features for ASR

- Equally-spaced frequency bands — but human hearing less sensitive at higher frequencies (above $\sim 1000\text{Hz}$)
- The estimated power spectrum contains harmonics of F_0 , which makes it difficult to estimate the envelope of the spectrum



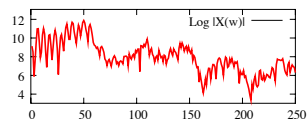
- Frequency bins of STFT are highly correlated each other, i.e. power spectrum representation is highly redundant

Cepstral Analysis

- Source-Filter model of speech production
 - **Source:** Vocal cord vibrations create a glottal source waveform
 - **Filter:** Source waveform is passed through the vocal tract: position of tongue, jaw, etc. give it a particular shape and hence a particular filtering characteristic
- Source characteristics (F_0 , dynamics of glottal pulse) do not help to discriminate between phones
- The filter specifies the position of the articulators
- ... and hence is directly related to phone discrimination
- Cepstral analysis enables us to separate source and filter

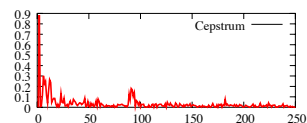
Cepstral Analysis

Split power spectrum into spectral envelope and F_0 harmonics.



Log Spectrum (freq domain)

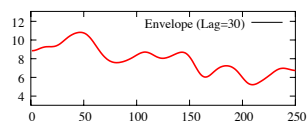
↓ Inverse Fourier Transform



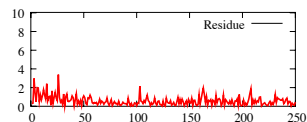
Cepstrum (time domain) (quefrequency)

↓ Liftering to get low/high part
(lifter: filter used in cepstral domain)

↓ Fourier Transform



Smoothed-spectrum (freq. domain)
[low-part of cepstrum]



Log spectrum
[high-part of cepstrum]

The Cepstrum

- Cepstrum obtained by applying inverse DFT to log magnitude spectrum (may be mel-scaled)
- Cepstrum is time-domain (we talk about quefrequency)
- Inverse DFT:

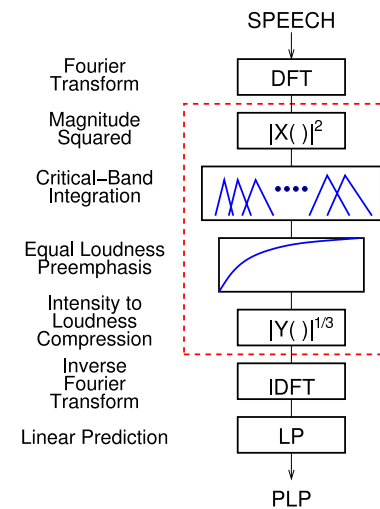
$$y_t[k] = \sum_{m=1}^M \log(|Y_t(m)|) \cos\left(k(m - 0.5)\frac{\pi}{M}\right) \quad K = 0, \dots, J$$

- Since log power spectrum is real and symmetric the inverse DFT is equivalent to a discrete cosine transform

MFCCs

- Smoothed spectrum: transform to cepstral domain, truncate, transform back to spectral domain
- Mel-frequency cepstral coefficients (MFCCs): use the cepstral coefficients directly
 - Widely used as acoustic features in HMM-based ASR
 - First 12 MFCCs are often used as the feature vector (removes F0 information)
 - Less correlated than spectral features — easier to model than spectral features
 - Very compact representation — 12 features describe a 20ms frame of data
 - For standard HMM-based systems, MFCCs result in better ASR performance than filter bank or spectrogram features
 - MFCCs are not robust against noise

PLP — Perceptual Linear Prediction



- PLP (Hermansky, JASA 1990)
- Uses equal loudness pre-emphasis and cube-root compression (motivated by perceptual results) rather than log compression
- Uses linear predictive auto-regressive modelling to obtain cepstral coefficients
- PLP has been shown to lead to
 - slightly better ASR accuracy
 - slightly better noise robustness
 compared with MFCCs

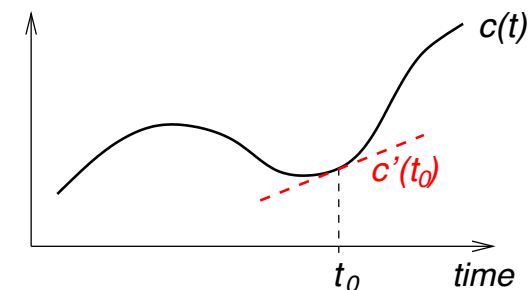
Dynamic features

- Speech is not constant frame-to-frame, so we can add features to do with how the cepstral coefficients change over time
- Δ^* , Δ^2* are delta features (dynamic features / time derivatives)
- Simple calculation of delta features $d(t)$ at time t for cepstral feature $c(t)$:

$$d(t) = \frac{c(t+1) - c(t-1)}{2}$$

- More sophisticated approach estimates the temporal derivative by using regression to estimate the slope (typically using 4 frames each side)
- “Standard” ASR features are 39 dimensions:
 - 12 MFCCs, and energy
 - 12 Δ MFCCs, Δ energy
 - 12 Δ^2 MFCCs, Δ^2 energy

Estimating dynamic features



Feature Transforms

- Orthogonal transformation (orthogonal bases)
 - **DCT** (discrete cosine transform)
 - **PCA** (principal component analysis)
- Transformation based on the bases that maximises the separability between classes.
 - **LDA** (linear discriminant analysis) / Fisher's linear discriminant
 - **HLDA** (heteroscedastic linear discriminant analysis)

Summary: Speech Signal Analysis for ASR

- Good characteristics of ASR features
- MFCCs - mel frequency cepstral coefficients
 - Short-time DFT analysis
 - Mel filter bank
 - Log magnitude squared
 - Inverse DFT (DCT)
 - Use first few (12) coefficients
- Delta features
- 39-dimension feature vector:
MFCC-12 + energy; + Deltas; + Delta-Deltas