# 9T1: Spectral-based audio features
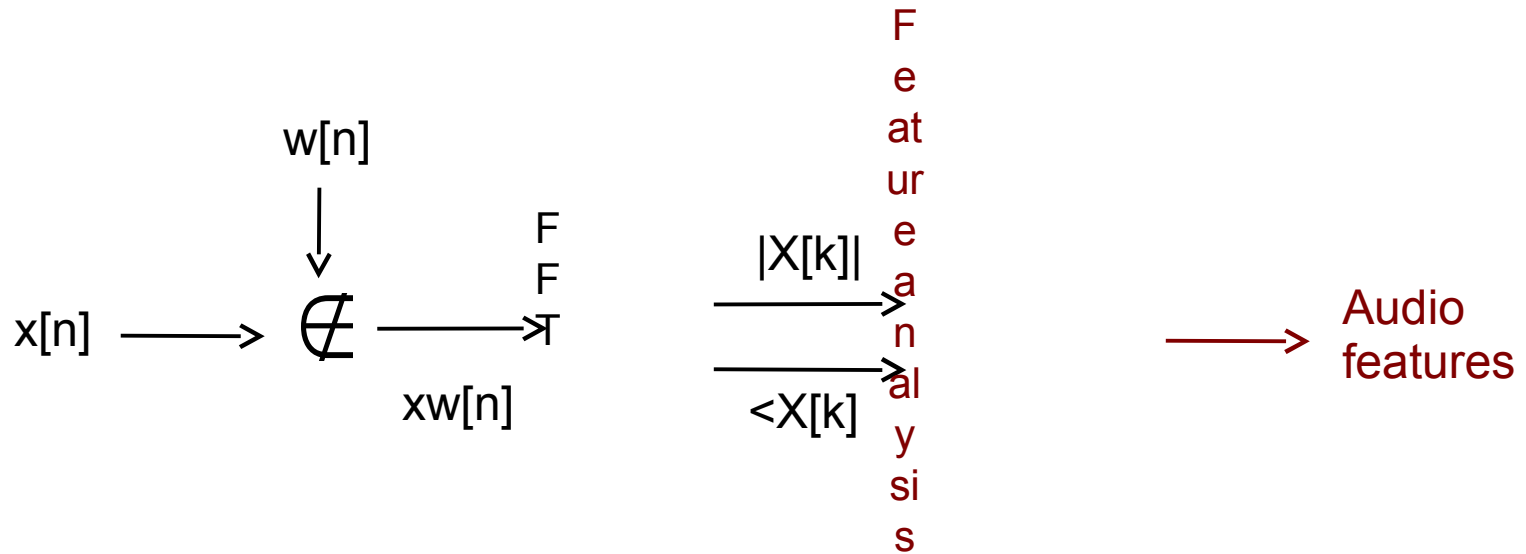
*Xavier Serra*

rsitat Pompeu Fabra, Barcelona

# Index

- Introduction: audio features

- Single-frame spectral features

- Multiple-frames spectral features

# Audio features

x[n] $\longrightarrow$

w[n] $\downarrow$

$\not\in$

xw[n] $\longrightarrow$ FFT

|X[k]| $\longrightarrow$

<X[k] $\longrightarrow$

Feature analysis

$\longrightarrow$ Audio features

# Essentia descriptors

- **Spectral descriptors:** `BarkBands, MelBands, ERBBands, MFCC, GFCC, LPC, HFC, SpectralContrast, Inharmonicity and Dissonance, ...`

- **Time-domain descriptors:** `EffectiveDuration, ZCR, Loudness, ...`

- **Tonal descriptors:** `PitchSalienceFunction, PitchYinFFT, HPCP, TuningFrequency, Key, ChordsDetection, ...`

- **Rhythm descriptors:** `BeatTrackerDegara, BeatTrackerMultiFeature, BpMHistogramDescriptors, NoveltyCurve, OnsetDetection, Onsets, ...`

- **SFX descriptors:** `LogAttackTime, MaxToTotal, MinToTotal, TCToTotal,...`

- **High-level descriptors:** `Danceability, DynamicComplexity, FadeDetection, SBic, ...`

# Single-frame spectral features

- Energy, RMS, Loudness

- Spectral centroid

- Mel-frequency cepstral coefficients (MFCC)

- Pitch salience

- Chroma (Harmonic pitch class profile, HPCP)
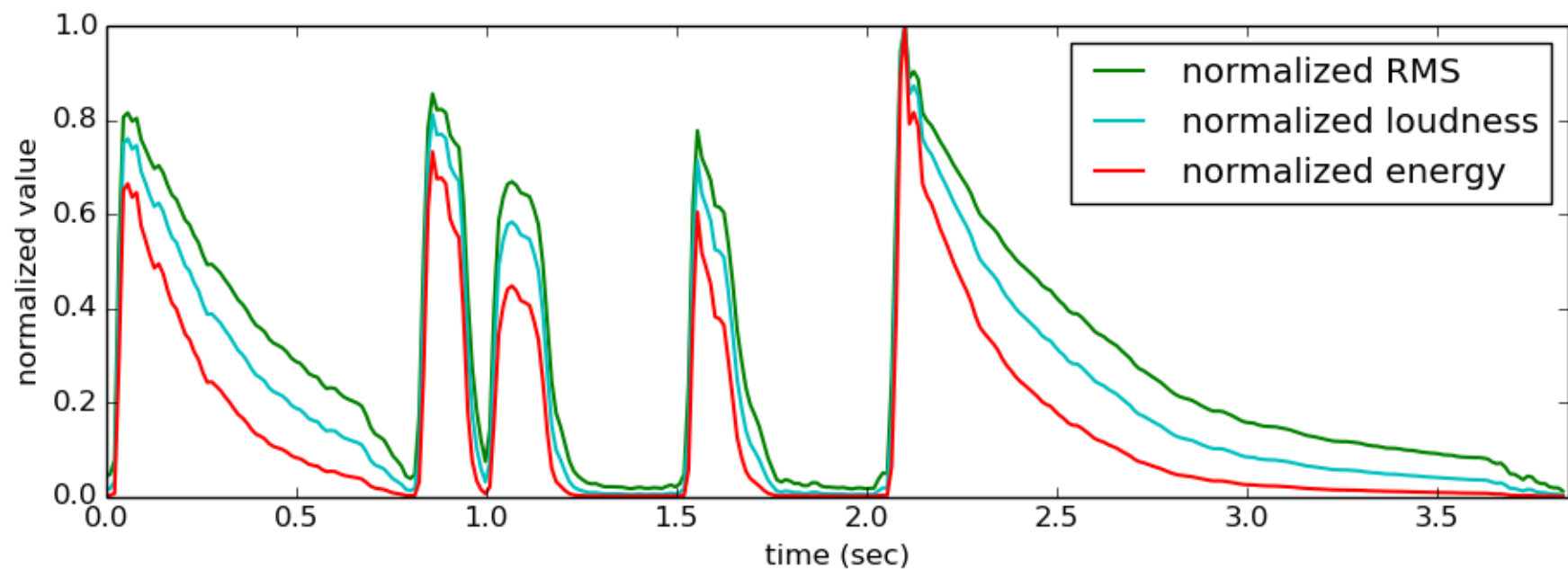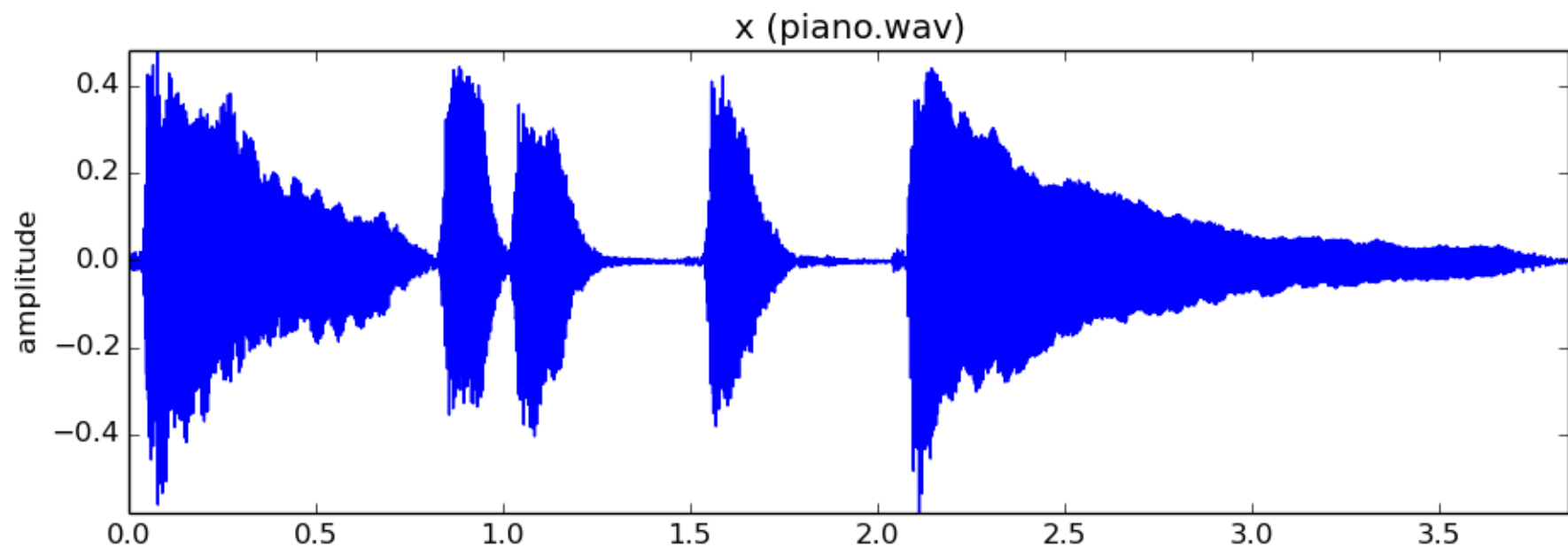
# Energy, RMS, Loudness

Energy:

$$energy_l = \sum_{k=0}^{N-1} \left( X_l[k] \right)^2$$

Root mean square:

$$RMS_l = \sqrt{\frac{1}{N^2} \sum_{k=0}^{N-1} \left( X_l[k] \right)^2}$$

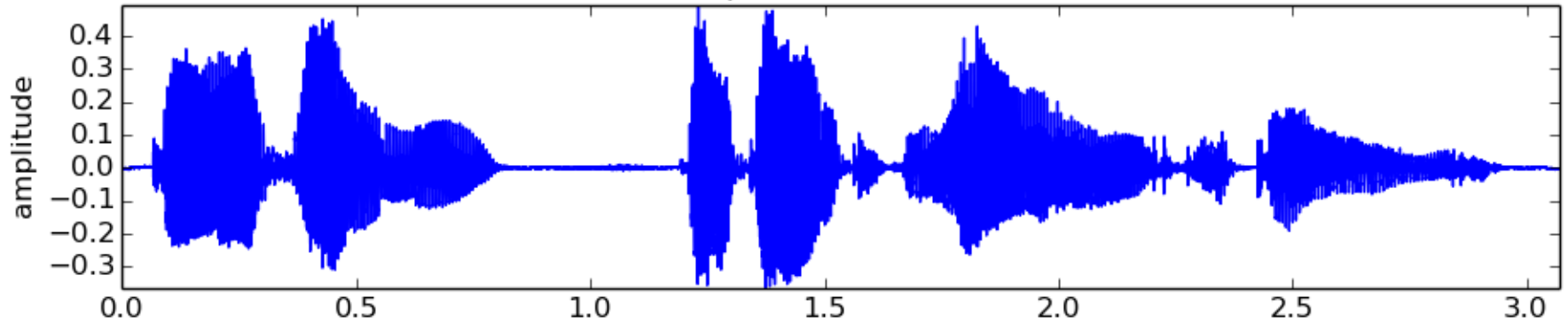Steven's power law:

$$loudness_l = \left( \sum_{k=0}^{N-1} \left( X_l[k] \right)^2 \right)^{0.67}$$

# Spectral centroid

$$centroid_l = \frac{\sum_{k=0}^{N/2} k \, |X_l[k]|}{\sum_{k=0}^{N/2} |X_l[k]|}$$



x (speech-male.wav)



spectral centroid

# Mel frequency cepstral coefficients

i is the number of output cepstral coefficients (goes up to H), H is the number of filters

INDEPENDENT FROM PITCH and LOUDNESS, it is a "spectrum shape feature"

$$mfcc_l = DCT\left(\log_{10}\left(\sum_{k=0}^{N/2} |X_l[k]||H_i[k]|\right)\right)$$

where

$|X[k]|$ is the positive magnitude spectrum
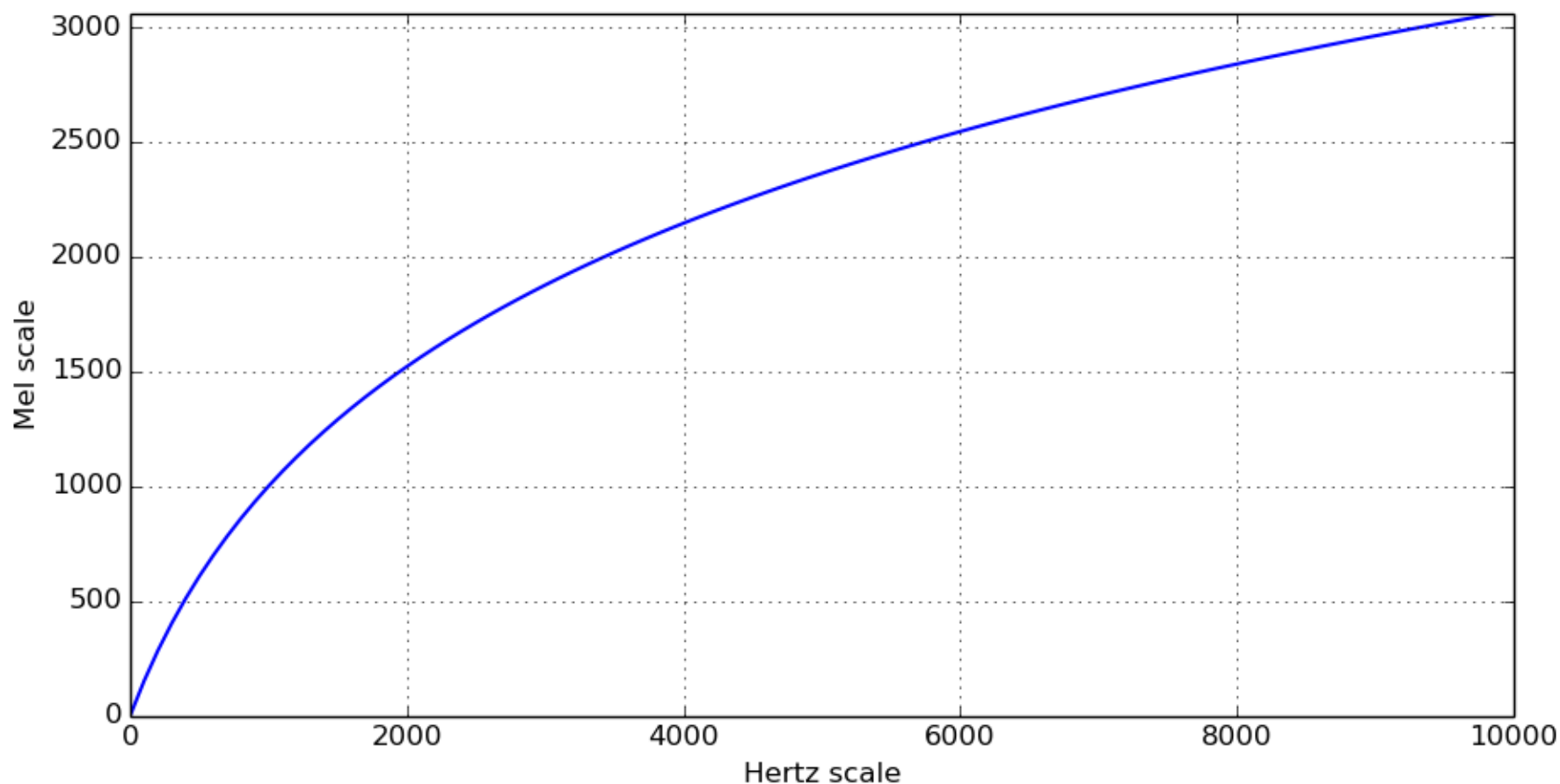
$H_i[k]$ is the mel scale filter bank for each filter i

$$DCT[m](\text{Discrete Cosine Transform}) = \sum_{n=0}^{N-1} f[n]\cos\left(\frac{\pi}{N}\left(n+\frac{1}{2}\right)m\right)$$

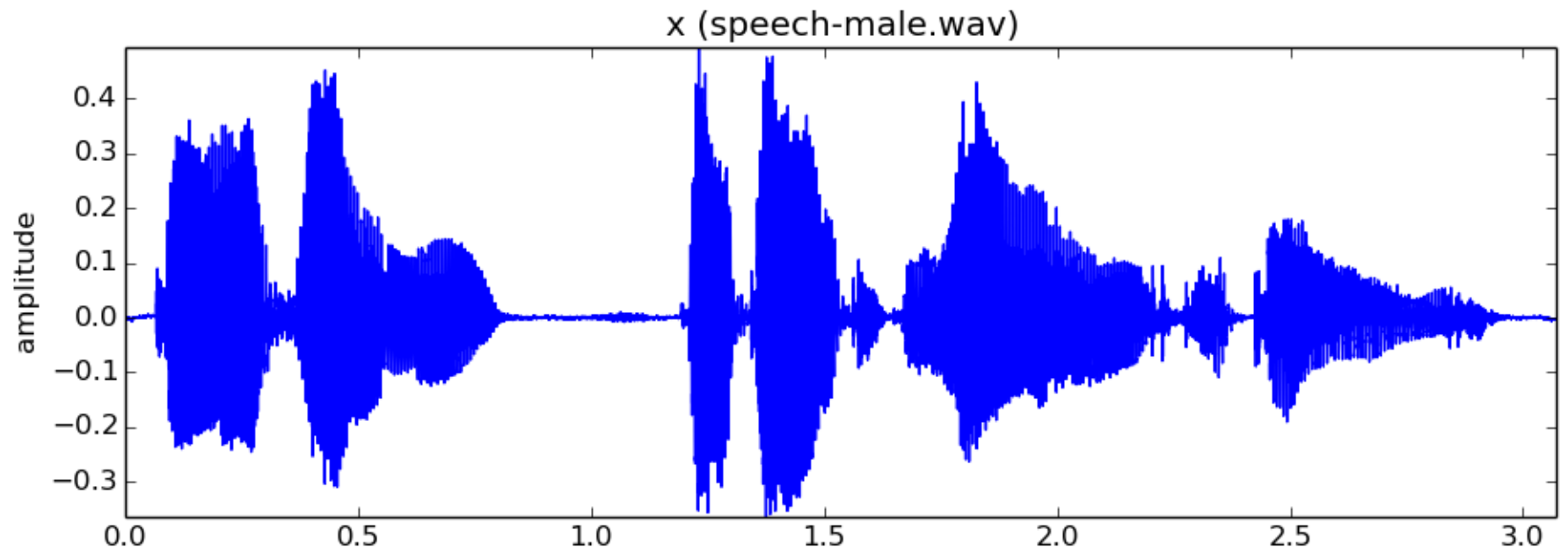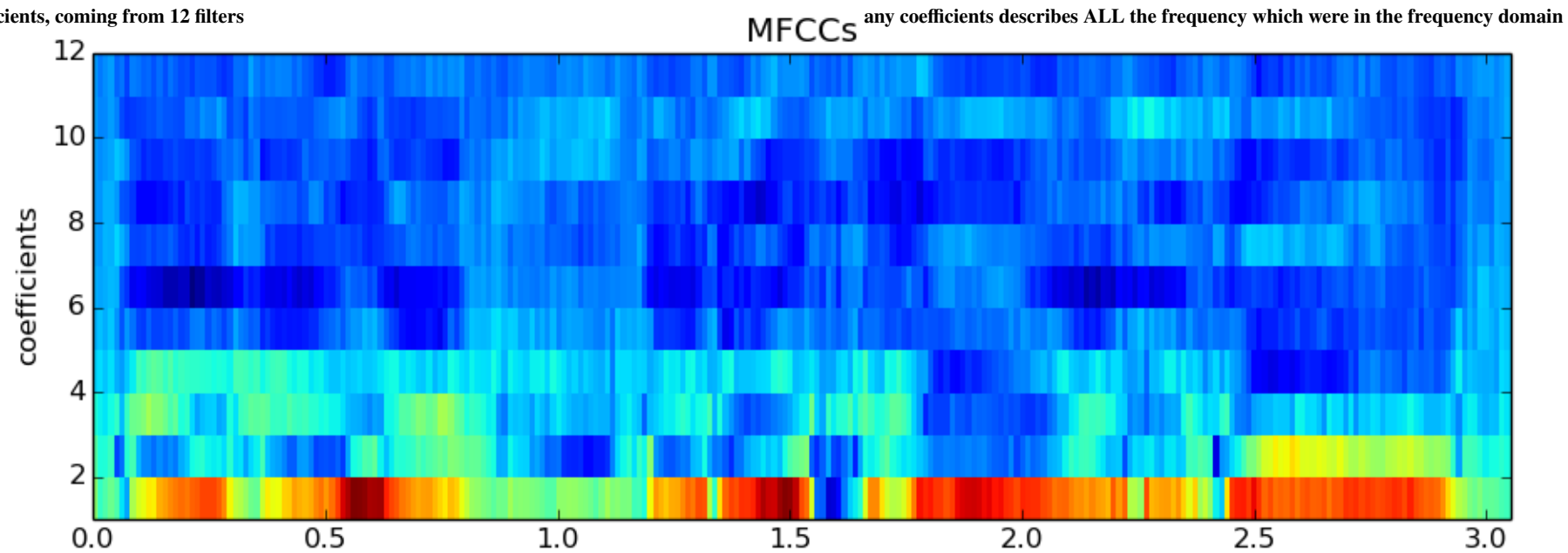"sampled version" of the original spectrum

linear magnitude spectrum

$|X_l[k]|$ $\longrightarrow$ H_i_[k_] $\longrightarrow$ L o g _ 1 0 $\longrightarrow$ D C T $\longrightarrow$ mfcc

apply some filters ( h )

The DCT (discrete cosine transform) of a DFT (the DCT of a frequency spectrum) compares shapes of sine waves -with different period- against the shape of the magnitude spectrum, and calculates the similarity between them

(usually) 12 coefficients in the cepstral domain

The DFT of a spectrum means comparing sine waves with different periods (0, 1, 2, ecc. times the number of frequency bins in the spectrum) against the SHAPE OF THE SPECTRUM

# MFCC: Mel scale

Perceptually relevant logarithmic scale

$$mel = 2595 \cdot \log_{10}\left(1 + \frac{f}{700}\right)$$

x (speech-male.wav)

**12 coefficients, coming from 12 filters** **any coefficients describes ALL the frequency which were in the frequency domain**

MFCCs

**coefficient n. 0 is an offset, DC, is the AMPLITUDE basically**

# Pitch salience

$\cup X_l[k] \cup$ → Peak detection → $A_p f_p$ → Pitch salience → $S_l[b]$

$$S[b] = \sum_{h=1}^{H} \sum_{p=1}^{P} e(A_p) g(b,h,f_p)(A_p)^{\beta}$$

where

$S[b] = $ salience at bin frequency b (b expressed in cent scale)
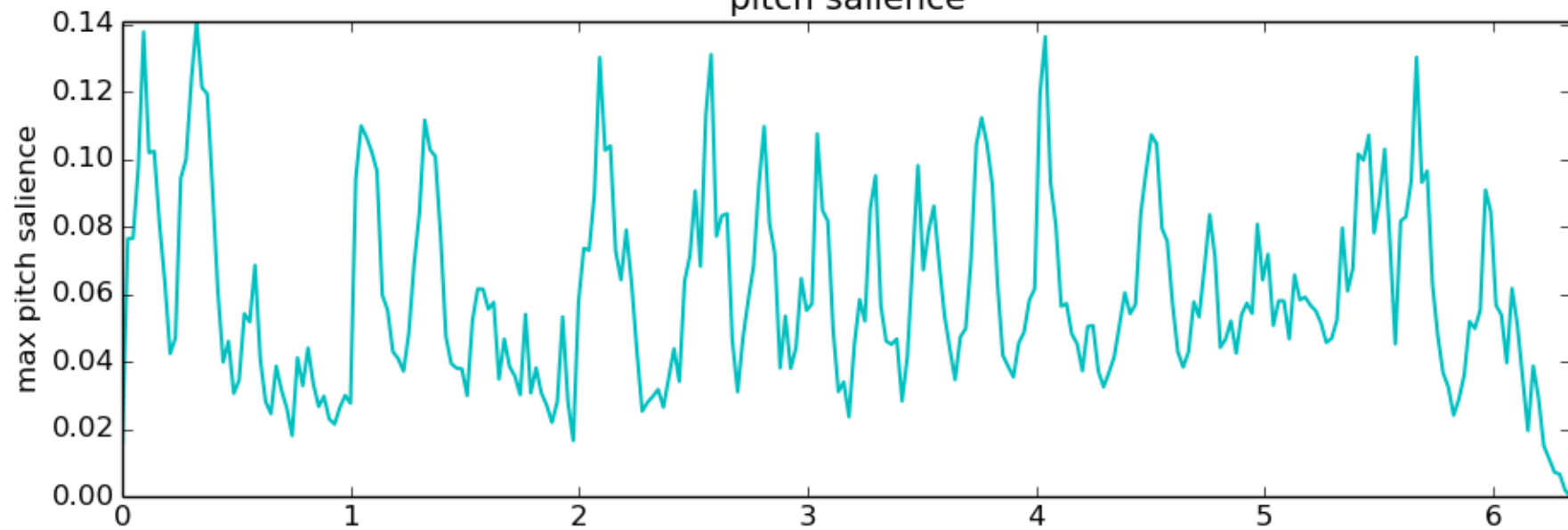
$e() = $ magnitude threshold function

$g() = $ weighting function applied to peak p

$\beta = $ magnitude compression value

# Chroma (Harmonic Pitch Class Profile)

$$hpcp[k] = \sum_{p=1}^{P} w(k, f_p) A_p^2$$

where

$A_p$ = amplitude of spectral peak $p$
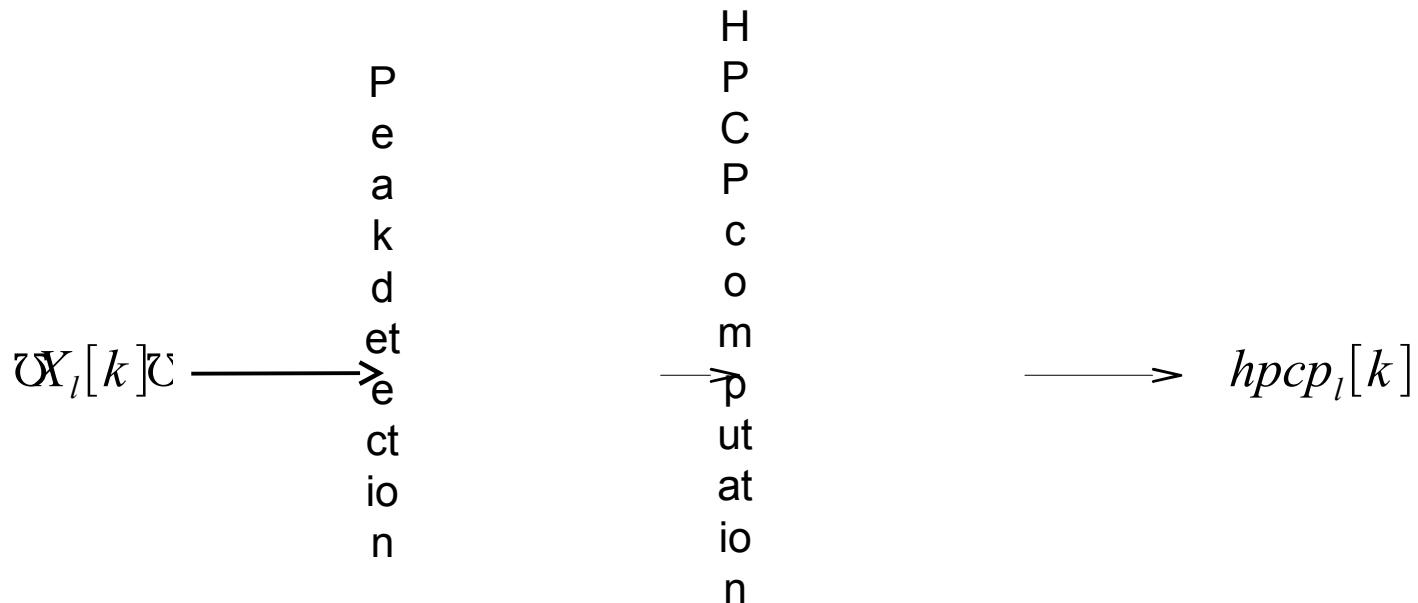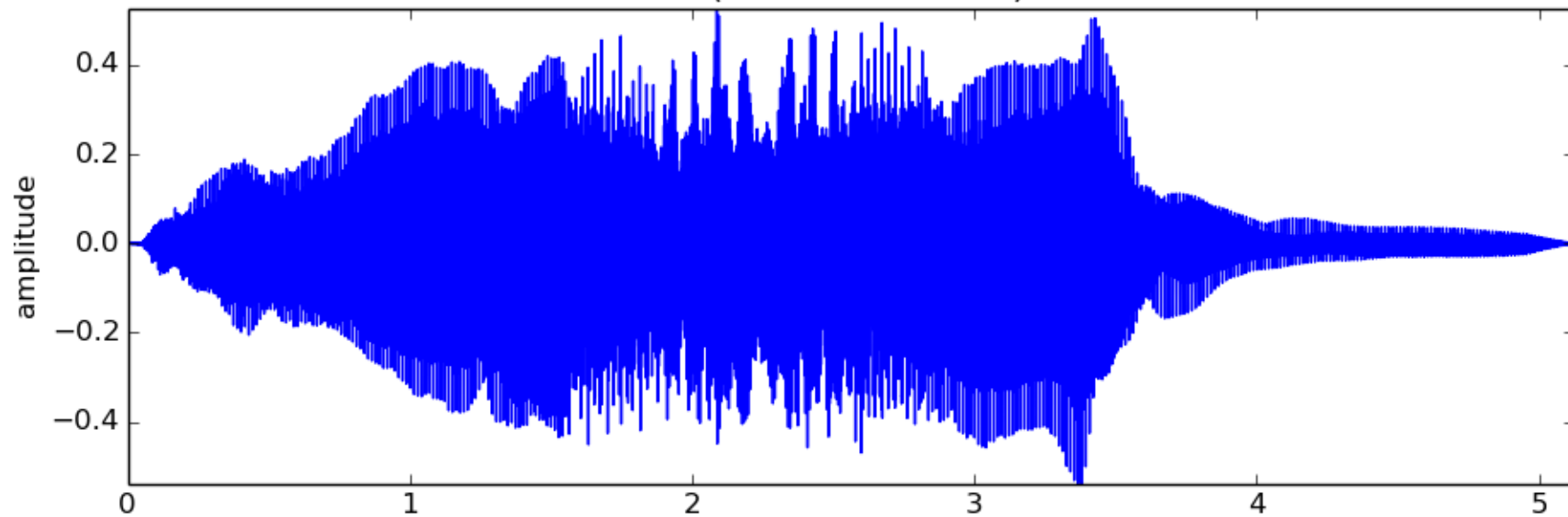
$w(k, f_p)$ = weight of the peak frequency f_p for bin k

$k$ = spectral bin locations of the chosen HPCP frequencies

$|X_l[k]|$ ⟶ Peak detection ⟶ HPCP computation ⟶ $hpcp_l[k]$

x (cello-double.wav)

HPCP

# Multiple-frames spectral features

- Event segmentation, onsets

- Predominant pitch

- Statistics of single-frame features

# Event segmentation, onsets

- Spectral flux (used in segmentation)

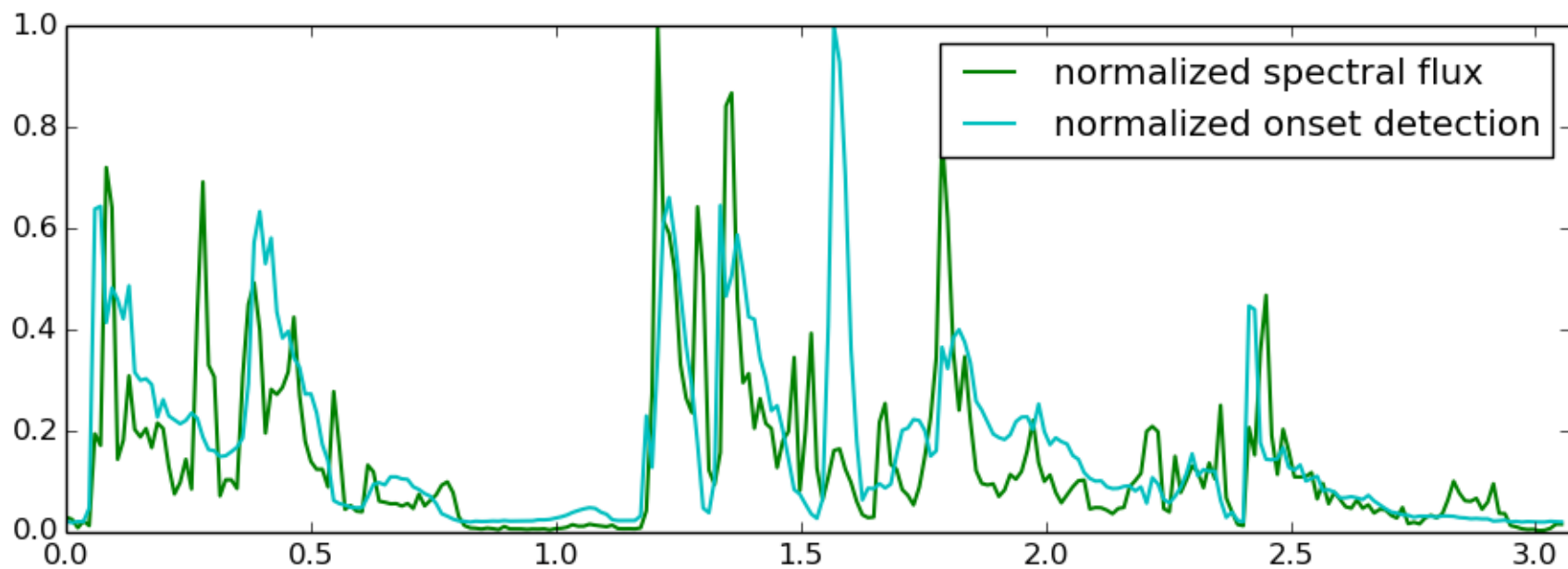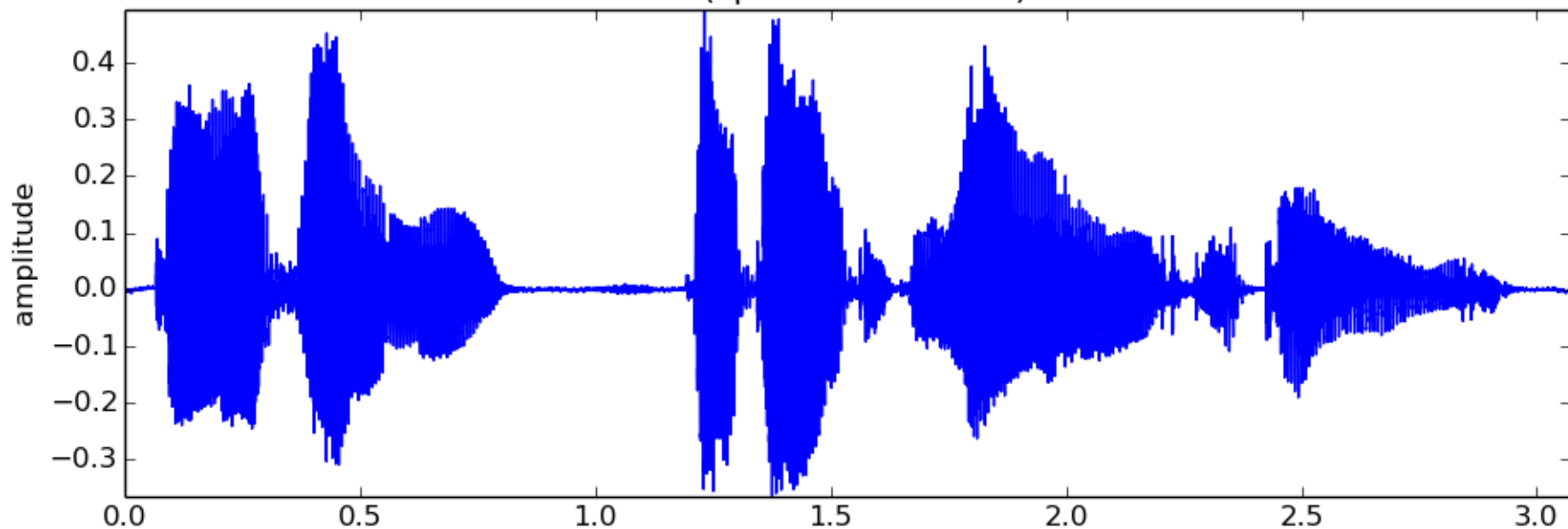$$SF_l = \sum_{k=0}^{N/2} H\left(\,\lvert X_l[k]\rvert - \lvert X_{(l-1)}[k]\rvert\,\right)$$

$$\text{where } H(x) = \frac{x + \lvert x \rvert}{2}$$

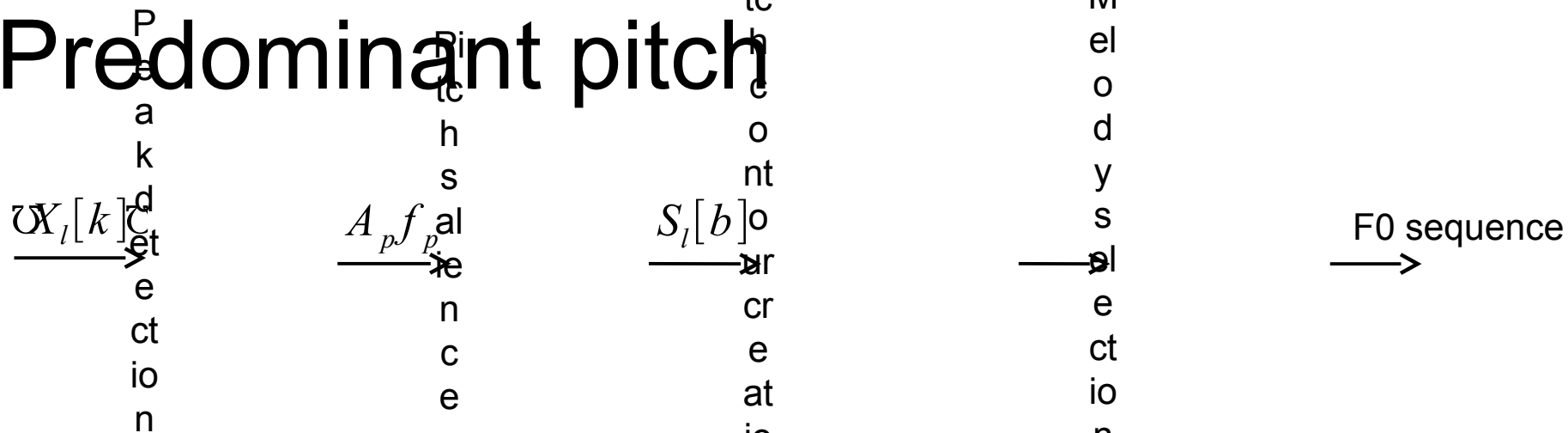- Onset detection based on high-frequency content

$$\text{Onset detection function} = HFC_l - HFC_{(l-1)}$$

$$\text{where } \quad HFC_l = \sum_{k=1}^{N/2} \lvert X_l[k]\rvert \cdot k^2$$

x (speech-male.wav)

— normalized spectral flux
— normalized onset detection

# Predominant pitch

$\underrightarrow{X_l[k]}$ **Peak detection** $\xrightarrow{A_p f_p}$ **Pitch salience** $\xrightarrow{S_l[b]}$ **Pitch contour creation** → **Melody selection** $\xrightarrow{\text{F0 sequence}}$



x (carnatic.wav)

amplitude



prominent melody

# Statistics of single frame features

- Arithmetic mean (first moment)

$$mean = \frac{1}{N} \sum_{i=0}^{N-1} y[i]$$

- Variance (second moment)

$$variance = \frac{1}{N} \sum_{i=0}^{N-1} (y[i] - mean)^2$$

- Skewness (third moment)

$$skewness = \frac{\frac{1}{N} \sum_{i=0}^{N-1} (y[i] - mean)^3}{[\frac{1}{N-1} \sum_{i=0}^{N-1} (y[i] - mean)^2]^{3/2}}$$

# References

- Essentia: http://essentia.upf.edu

- http://en.wikipedia.org/wiki/Spectral_centroid

- http://en.wikipedia.org/wiki/Mel-frequency_cepstrum

- http://en.wikipedia.org/wiki/Loudness

- http://en.wikipedia.org/wiki/Harmonic_pitch_class_profiles

- http://en.wikipedia.org/wiki/Onset_(audio)

- http://en.wikipedia.org/wiki/Moment_(mathematics)

- Slides released under CC Attribution-Noncommercial-Share Alike license and code under Affero GPL license; available from https://github.com/MTG/sms-tools

# **9T1:** Spectral-based audio features

*Xavier Serra*

rsitat Pompeu Fabra, Barcelona