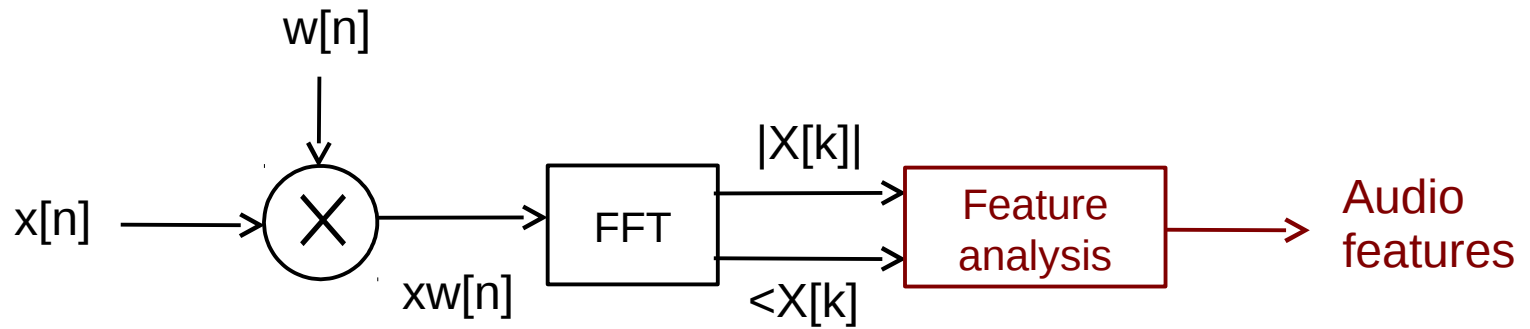# **9T1:** Spectral-based audio features

## *Xavier Serra*

Universitat Pompeu Fabra, Barcelona

# Index

- Introduction: audio features

- Single-frame spectral features

- Multiple-frames spectral features

# 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)
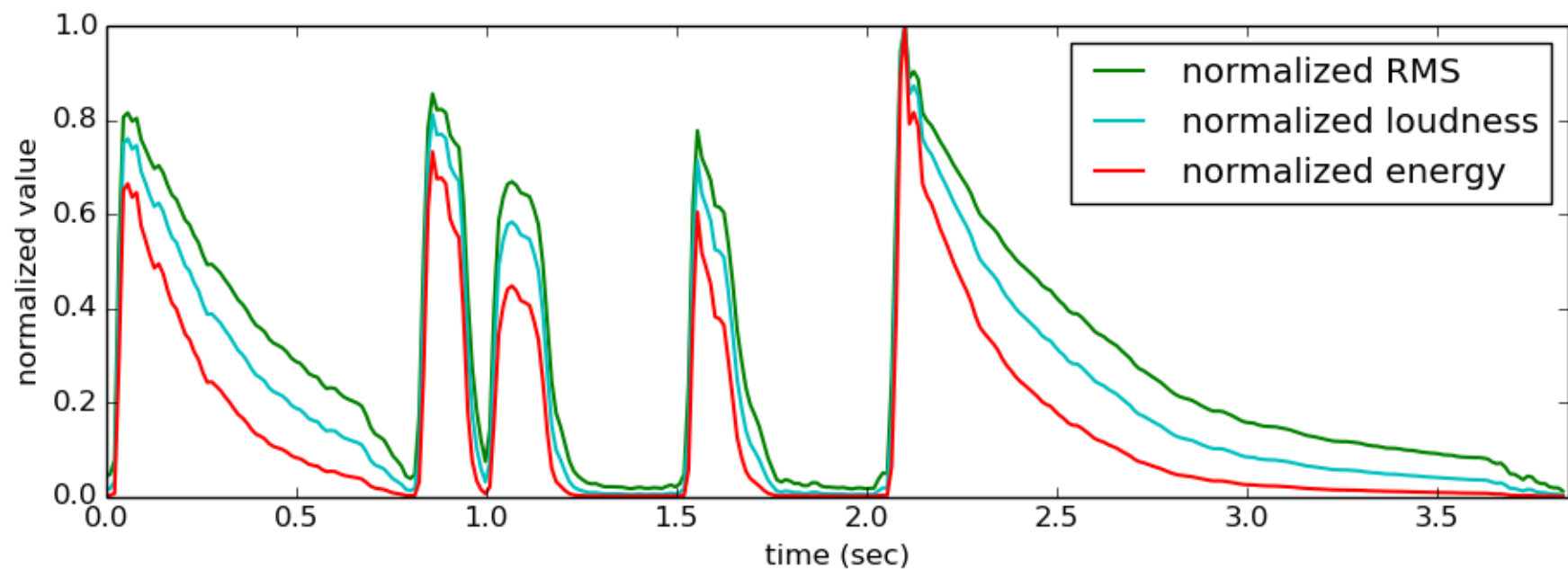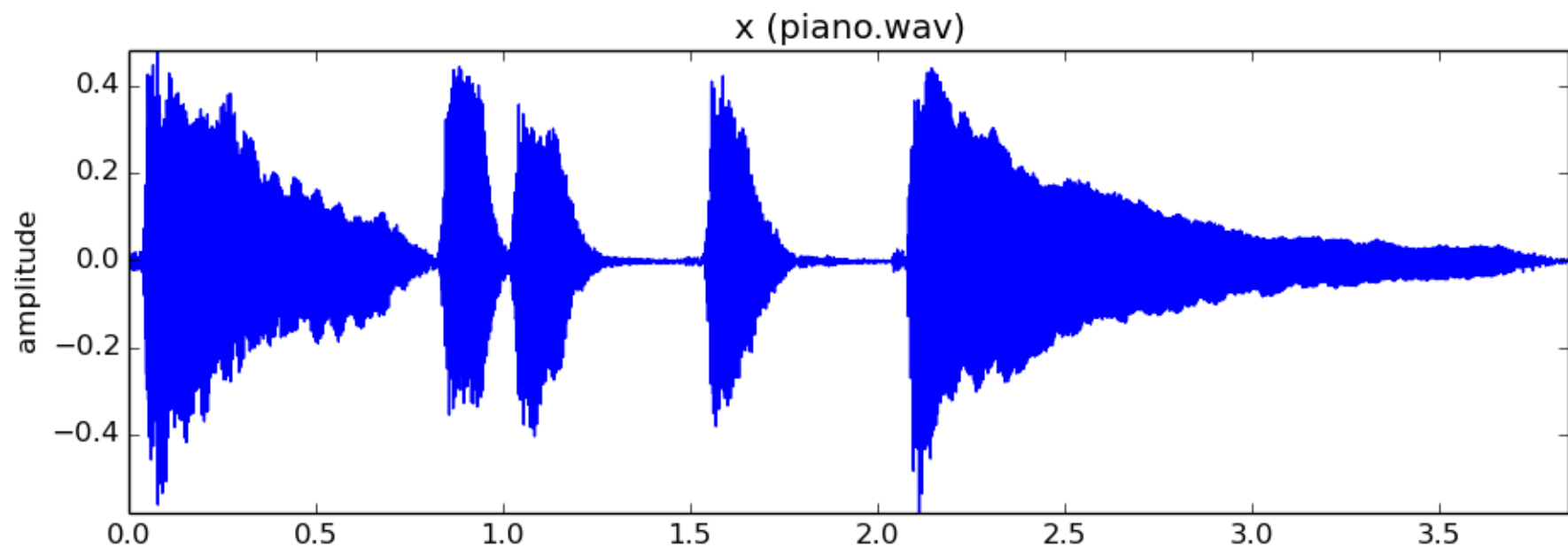
# Energy, RMS, Loudness

Energy:

$$energy_l = \sum_{k=0}^{N-1} |X_l[k]|^2$$

Root mean square:

$$RMS_l = \sqrt{\frac{1}{N^2} \sum_{k=0}^{N-1} |X_l[k]|^2}$$
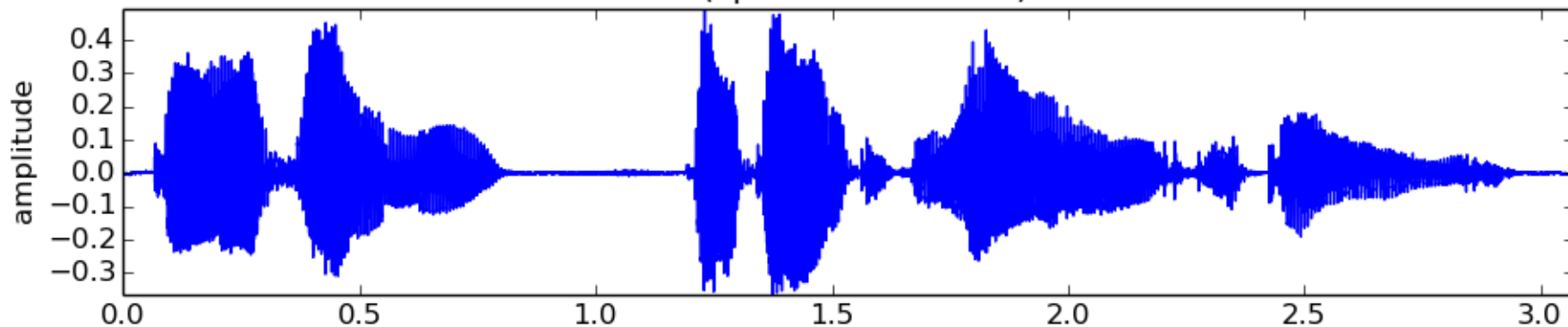
Steven's power law:

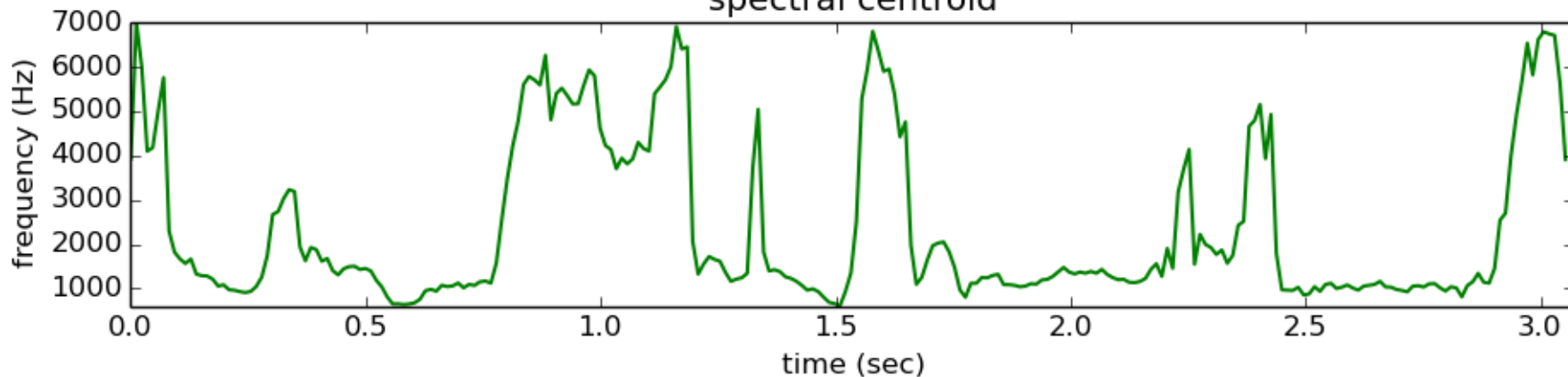$$loudness_l = \left( \sum_{k=0}^{N-1} |X_l[k]|^2 \right)^{0.67}$$

# Spectral centroid

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



x (speech-male.wav)

spectral centroid

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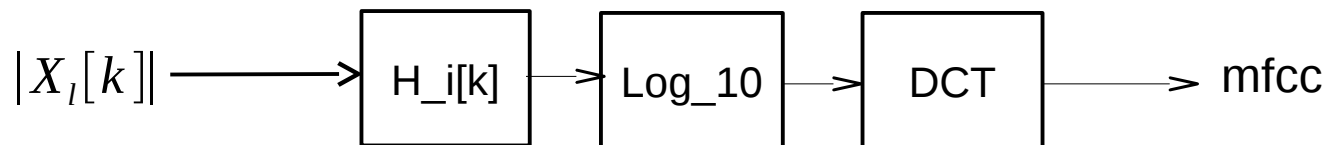
$$mfcc_l = DCT\left(\log_{10}\left(\sum_{k=0}^{N/2} |X_l[k]| |H_i[k]|\right)\right)$$
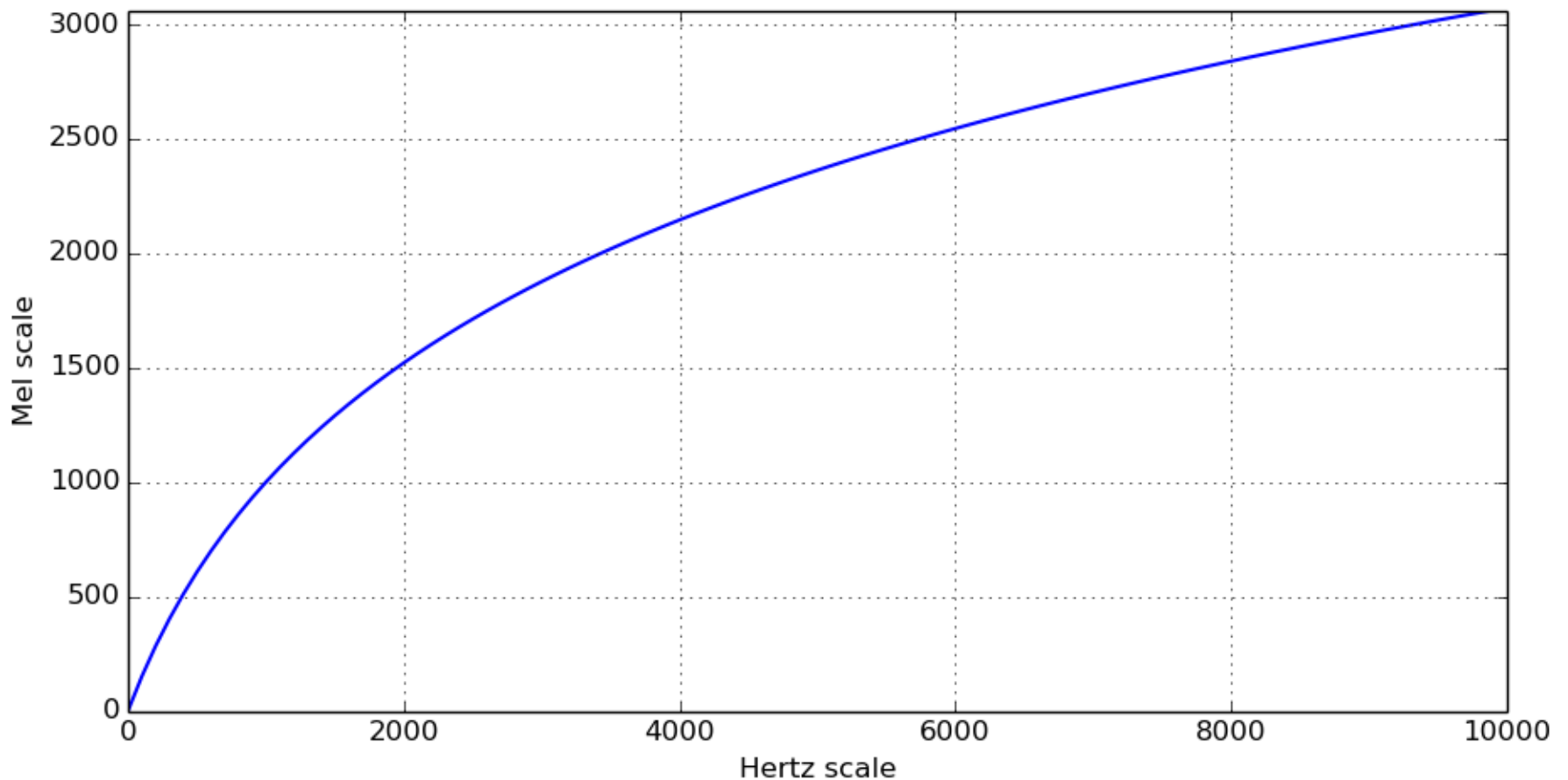
where

$|X[k]|$ is the positive magnitude spectrum

$H_i[k]$ is the mel scale filter bank for each filter i

$$DCT[m](\text{Discrete Cosine Transform}) = \sum_{n=0}^{N-1} f[n]\cos\left(\frac{\pi}{N}\left(n+\frac{1}{2}\right)m\right)$$
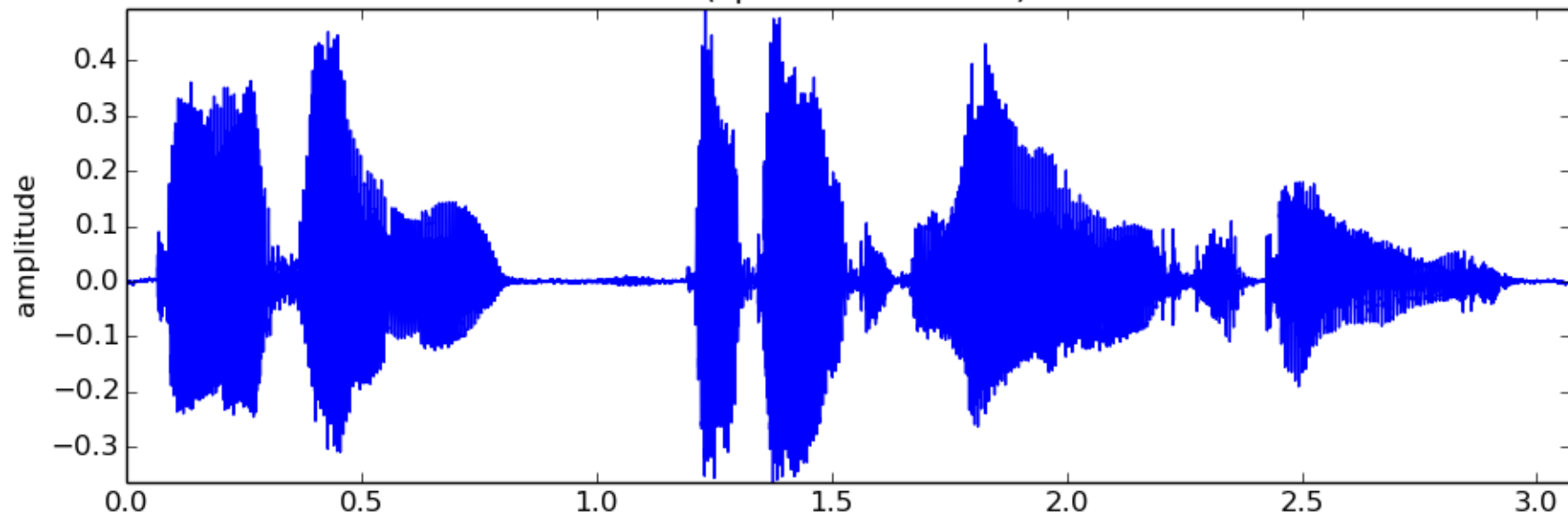
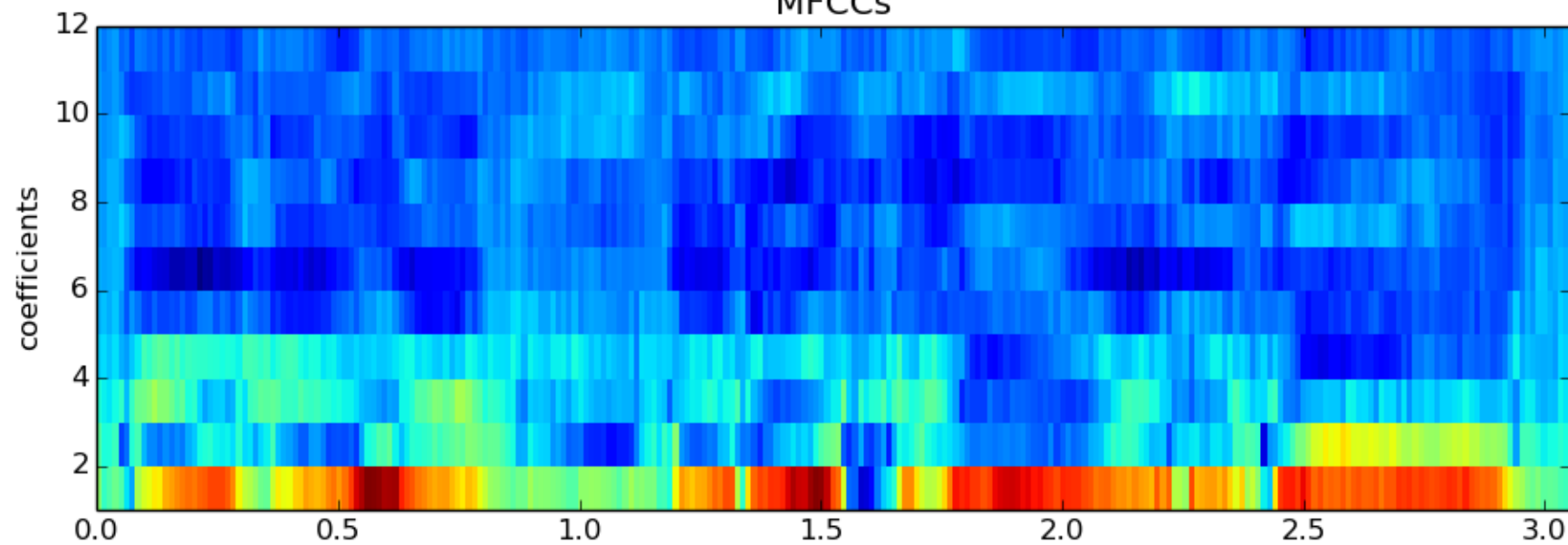$|X_l[k]| \longrightarrow$ H_i[k] $\rightarrow$ Log_10 $\rightarrow$ DCT $\longrightarrow$ mfcc

# MFCC: Mel scale

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

x (speech-male.wav)

MFCCs

# Pitch salience

$$|X_l[k]| \xrightarrow{\quad} \boxed{\begin{array}{c} \text{Peak} \\ \text{detection} \end{array}} \xrightarrow{A_p f_p} \boxed{\begin{array}{c} \text{Pitch} \\ \text{salience} \end{array}} \xrightarrow{S_l[b]}$$

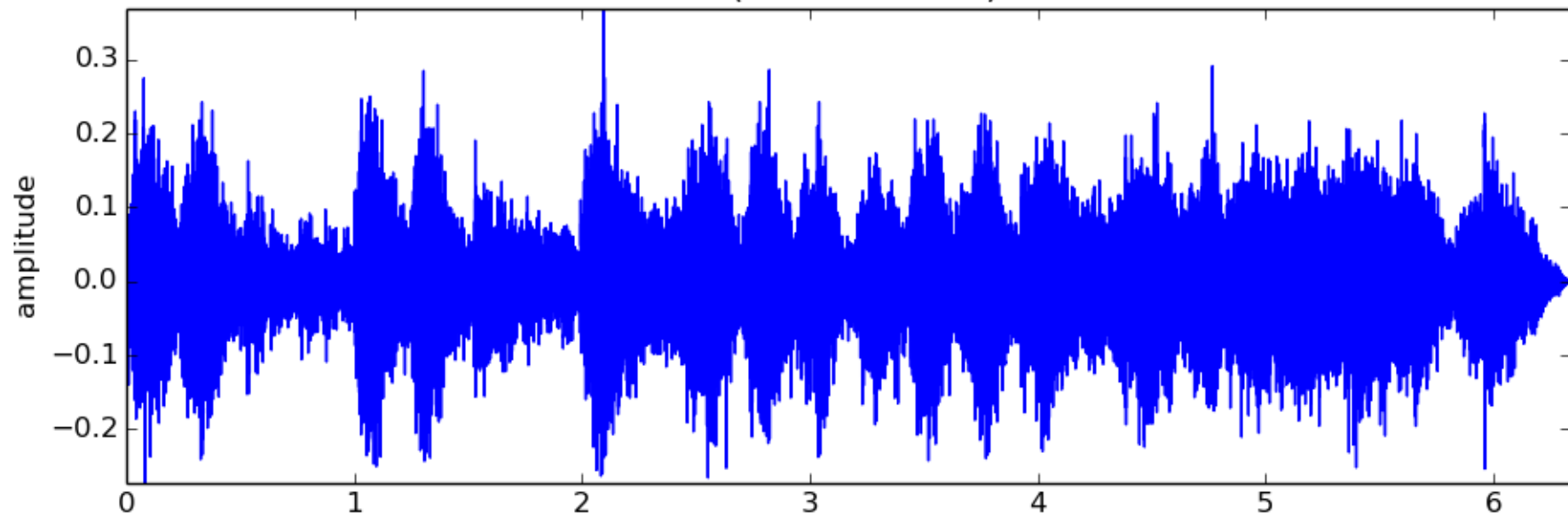$$S[b] = \sum_{h=1}^{H} \sum_{p=1}^{P} e(A_p) g(b, h, f_p)(A_p)^{\beta}$$

where

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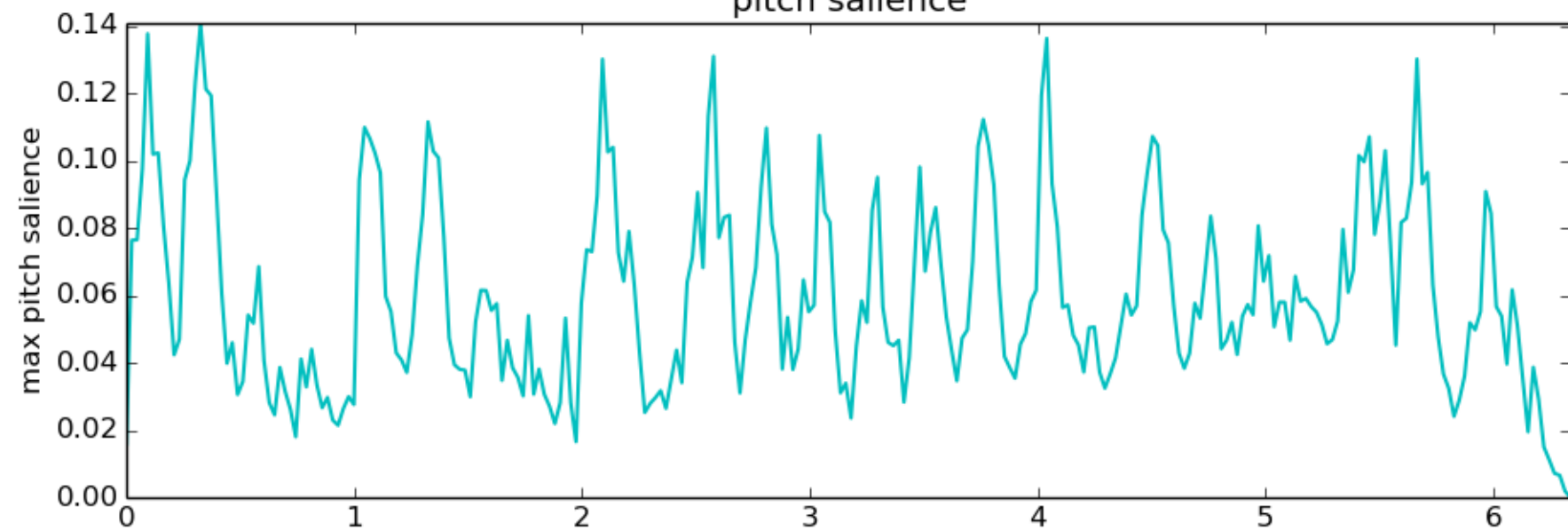
$e() =$ magnitude threshold function

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

$\beta =$ magnitude compression value
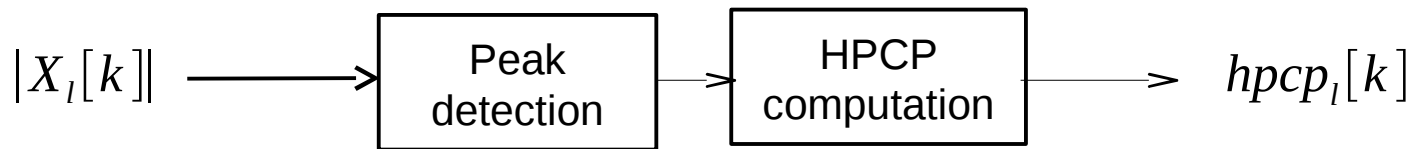
# Chroma (Harmonic Pitch Class Profile)

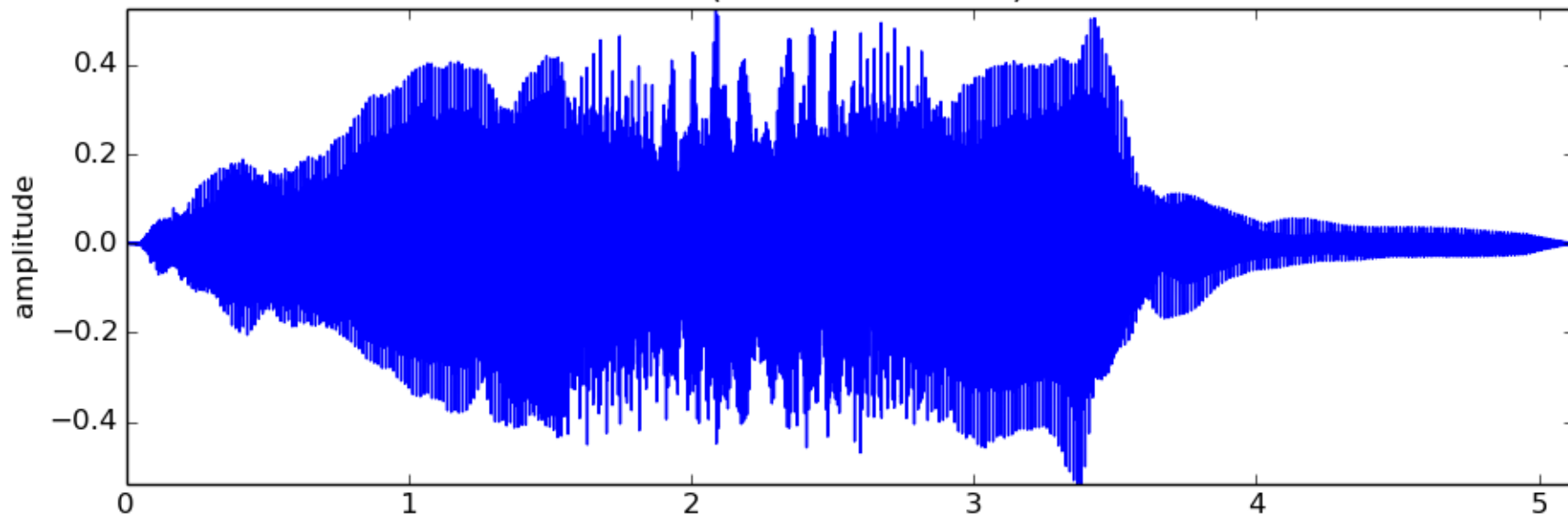$$hpcp[k] = \sum_{p=1}^{P} w(k, f_p) A_p^2$$

where

   $A_p =$ amplitude of spectral peak $p$

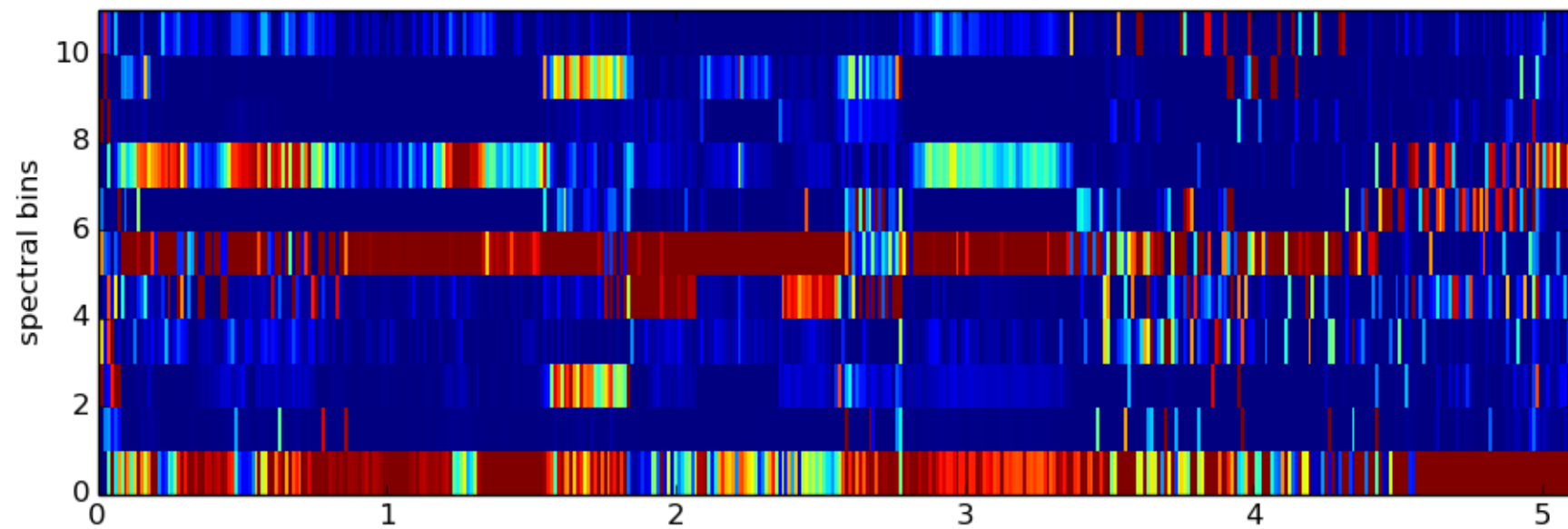   $w(k, f_p) =$ weight of the peak frequency f_p for bin k

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

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

# Multiple-frames spectral features

- Event segmentation, onsets

- Predominant pitch

- Statistics of single-frame features

# Event segmentation, onsets

- Spectral flux (used in segmentation)

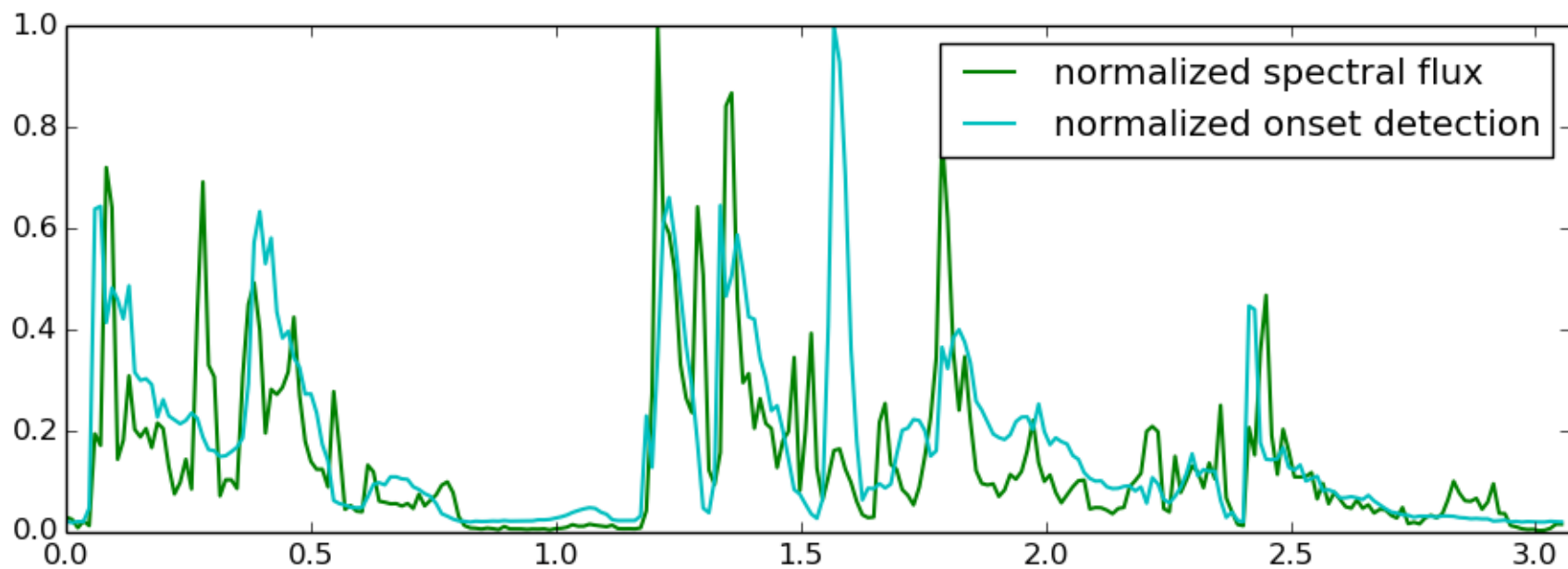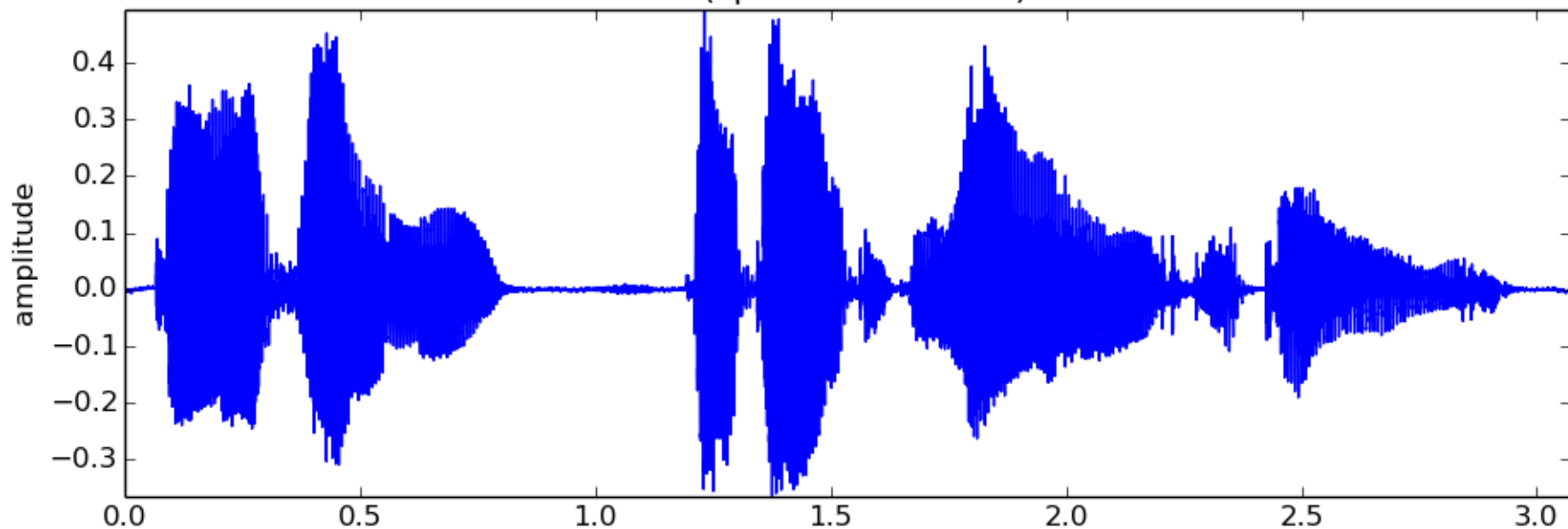$$SF_l = \sum_{k=0}^{N/2} H\left(\left|X_l[k]\right| - \left|X_{(l-1)}[k]\right|\right)$$

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

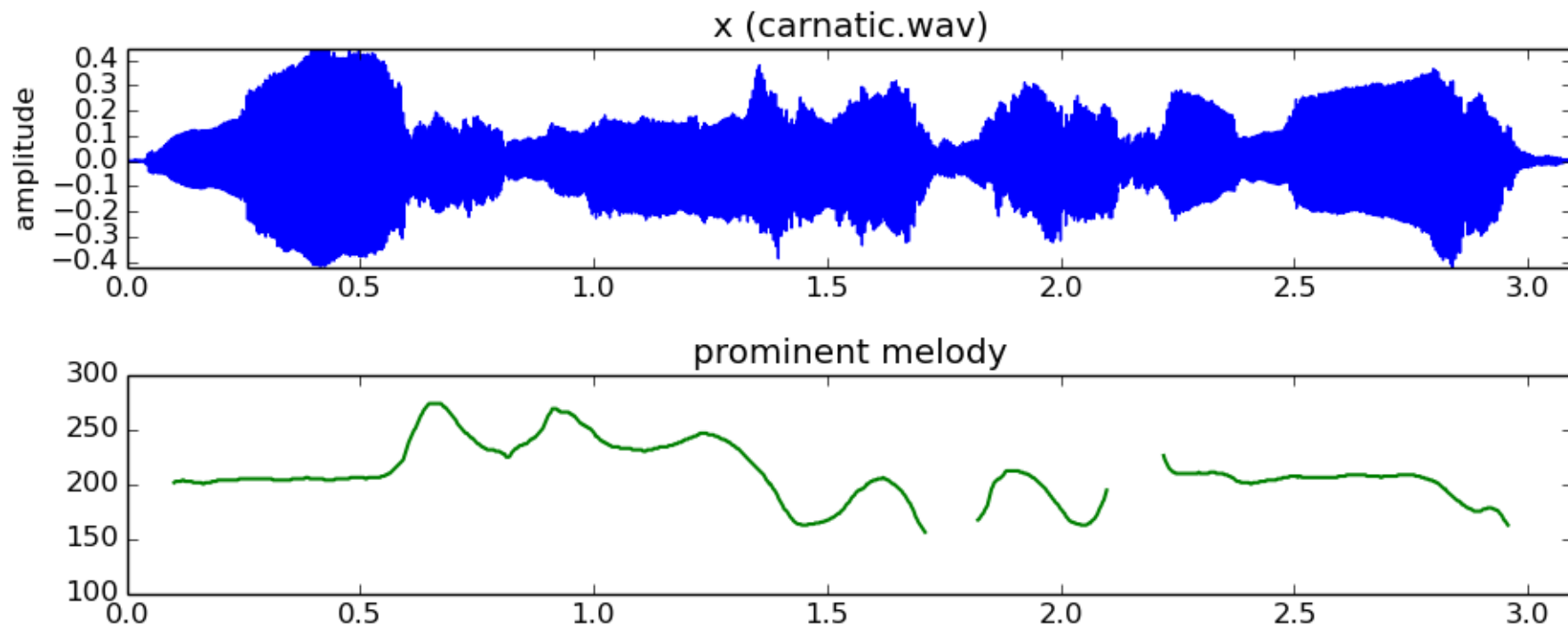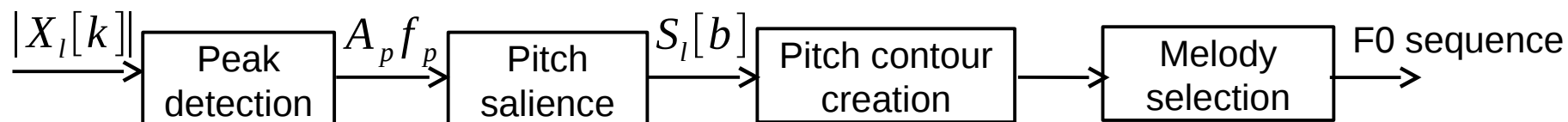- Onset detection based on high-frequency content

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

$$\text{where} \quad HFC_l = \sum_{k=1}^{N/2} \left|X_l[k]\right| k^2$$

x (speech-male.wav)

- normalized spectral flux
- normalized onset detection

# Predominant pitch

$|X_l[k]|$ → [ Peak detection ] → $A_p f_p$ → [ Pitch salience ] → $S_l[b]$ → [ Pitch contour creation ] → [ Melody selection ] → F0 sequence



x (carnatic.wav)

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*

Universitat Pompeu Fabra, Barcelona