# Spectral based sound description

## *Xavier Serra*

Music Technology Group

Universitat Pompeu Fabra, Barcelona

*http://mtg.upf.edu*

# Index

- Sinusoidal plus residual model features

- Spectral-based features in Essentia

- Features for sound/music description

- Features for instrument modeling

- Music description

# Sinusoidal+residual model features

- Instantaneous frequency and amplitude of partials

- Instantaneous spectrum of residual

- Instantaneous fundamental frequency

- Amplitude and spectral shape of sinusoidal component

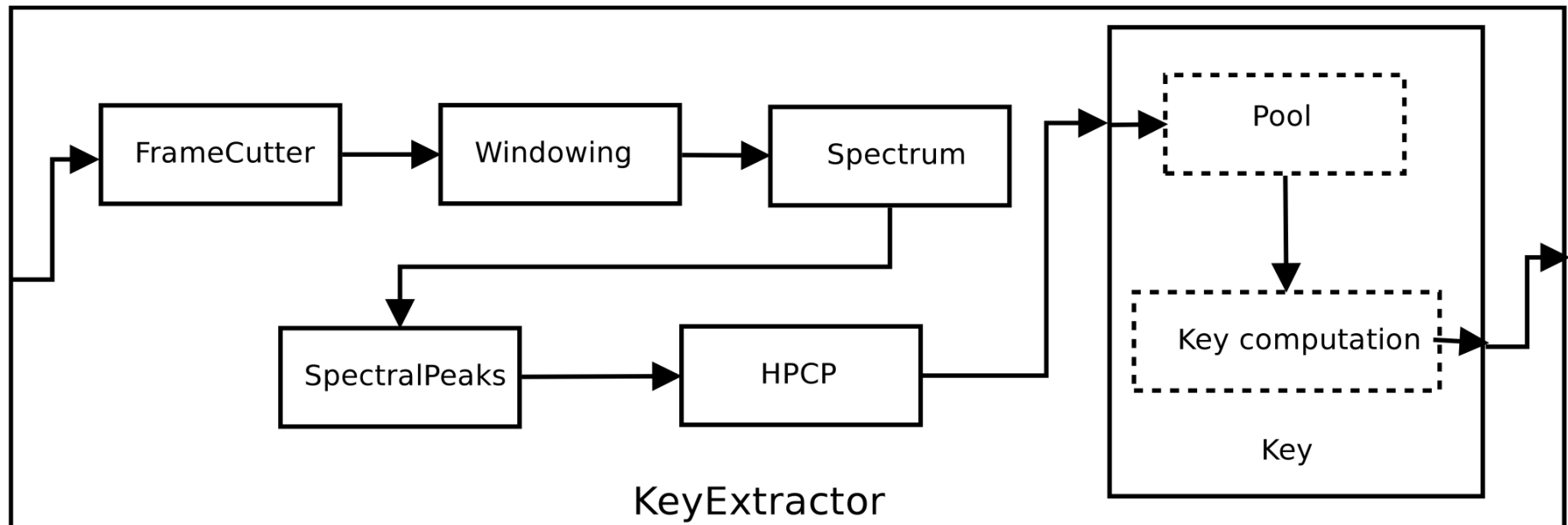- Amplitude and spectral shape of residual

# Essentia functionality

- audio file i/o; standard DSP building blocks; filters
- Descriptors such as
  - <u>spectral</u>: spectral shape, MFCC, Bark/Mel bands
  - <u>time-domain/rhythmic</u>: loudness, dynamics, onsets, beats, beats per minute, danceability
  - <u>tonal</u>: melody, pitch, chroma, key, scale, chords
  - <u>high-level</u>: segmentation, genres, mood (happy/sad), instrumentation (acoustic, electronic, timbre dark/bright, voice male/female)

# Spectral features in Essentia

- **BarkBands:** computes the Bark band energies.
- **MelBands:** computes the Mel band energies.
- **ERBBands:** computes the energies in bands spaced on an Equivalent Rectangular Bandwidth scale.
- **MFCC:** computes the Mel-frequency cepstral coefficients of a frame.
- **GFCC:** computes the gammatone feature cepstrum coefficients similar to MFCCs.
- **LPC:** computes the Linear Predictive Coding coefficients of a frame as well as the associated reflection coefficients.
- **HFC:** computes the High-Frequency Content measure.
- **SpectralContrast**: computes spectral contrast of a spectrum.
- **Inharmonicity and Dissonance:** both try to estimate whether an audio frame "sounds" harmonic or not.
- **SpectralWhitening:** whitens the input spectrum.
- **Panning:** computes the panorama distribution of a stereo audio frame.

# Essentia extractors

Executable extractors built by combining algorithms in a "data-flow" manner

```
import essentia
from essentia.standard import *
from pylab import *

loader = essentia.standard.MonoLoader(filename = 'oboe.wav')
audio = loader()

w = Windowing(type = 'hann')
spectrum = Spectrum()
mfcc = MFCC()

pool = essentia.Pool()

for frame in FrameGenerator(audio, frameSize = 1024, hopSize = 512):
    mfcc_bands, mfcc_coeffs = mfcc(spectrum(w(frame)))
    pool.add('lowlevel.mfcc', mfcc_coeffs)
    pool.add('lowlevel.mfcc_bands', mfcc_bands)

output = YamlOutput(filename = 'mfcc.sig')
output(pool)
```

# Sound similarity

# Prominent pitch detection



(Salamon, 2013)

# Peak tracking in polyphonic signals

# Onset detection

# Instrument model

- Spectral Shape Models (formants, average shape, ...)
- Phase Models (constant, formant-based, ...)
- Frequency Models (harmonic model, piano model)
- Vibrato Models
- Articulation Models (frequency, amplitude functions)
- Residual Models
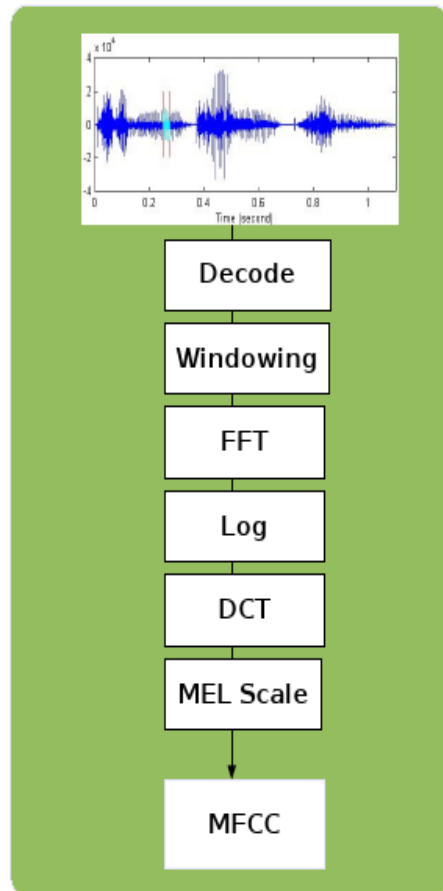- Brightness-Loudness model (amplitude versus spectral tilt)
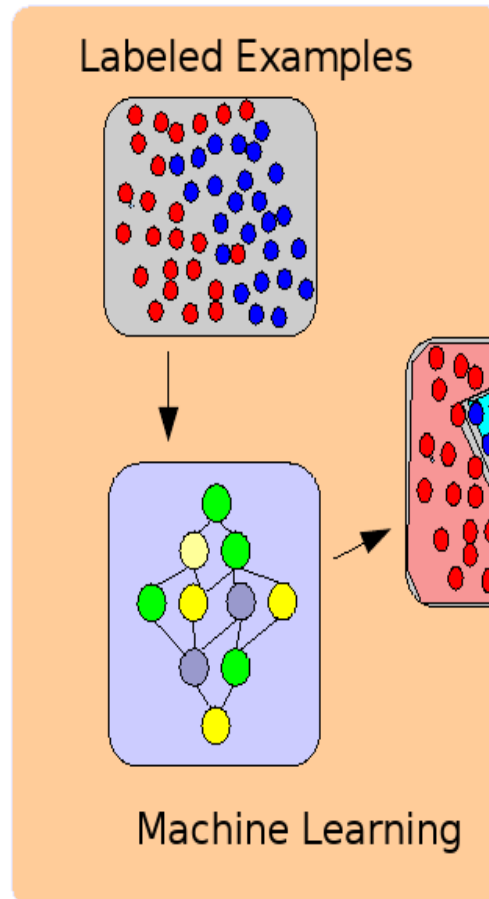
# Taxonomy of musical features

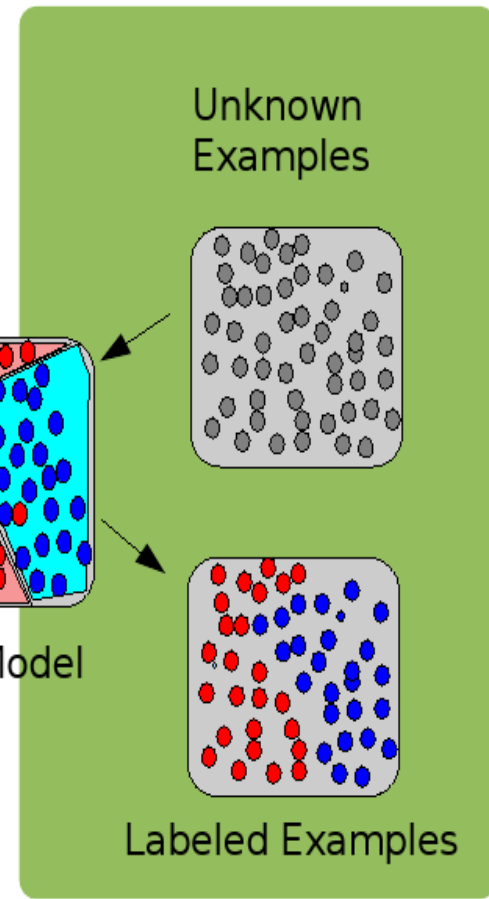| STRUCT | | CONCEPT LEVEL | | MUSICAL CONTENT FEATURES | | | | |
|---|---|---|---|---|---|---|---|---|
| **CONTEXTUAL** | global beyond 3 sec | **HIGH II** | **EXPRESSIVE** | cognition \| emotion \| affect = *syntactic+semantic concepts* | | | | |
| | | **HIGH I** | **FORMAL** | melody | harmony | rhythm | source | dynamics |
| | | | | key | tonality | rhythmic patterns | instrument | trajectory |
| | | | | profile | cadence | tempo | voice | articulation |
| | global < 3 sec | **MID** | **PERCEPTUAL** | successive intervallic pattern | simultane intervallic pattern | beat | spectral envelope | dynamic range |
| | | | | | | IOI | | sound level |
| | | | | pitch | | time | timbre | loudness |
| **NON-CONTEXTUAL** | local + spatial | **LOW II** | **SENSORIAL** | periodicity pitch | | note duration | roughness | neural energy |
| | | | | pitch deviations | | onset | spectral flux | peak |
| | local + temporal | **LOW I** | **PHYSICAL** | fundamental frequency | | offset | spectral centroid | |
| | | | | frequency | | duration | spectrum | intensity |

Lesaffre et alt., 2003
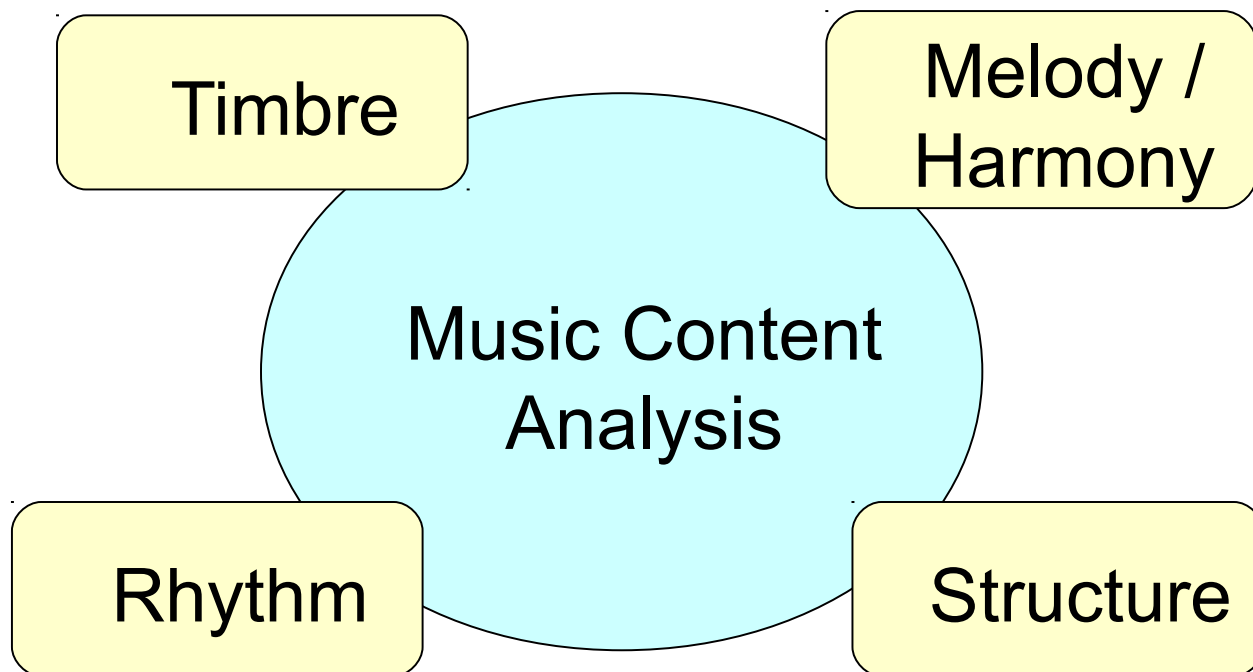
# Audio content classification

# Levels of description

- Low-level (signal-centered) descriptors: computed from the audio signal in a direct or derived (ex: spectral analysis) way: average energy, spectral centroid, MFCCs ….

- Mid-level (object-centered) descriptors: requiring an induction operation or data modeling: key, genre, instrument …

- High-level (user-centered) descriptors: requiring a user model: mood (ex: happy, sad), …
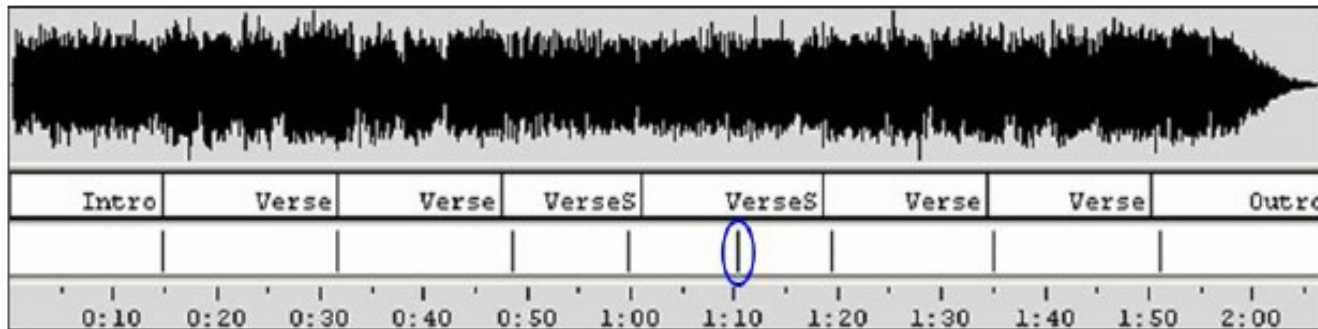
# Facets of music content

# Structure description

- Partitioning the sound stream into *homogeneous* regions

- Detecting special roles for the segmented regions: intro, verse, chorus, bridge,

- Other segments can also be identified: instrumental / singing; solo / ensemble; chords…
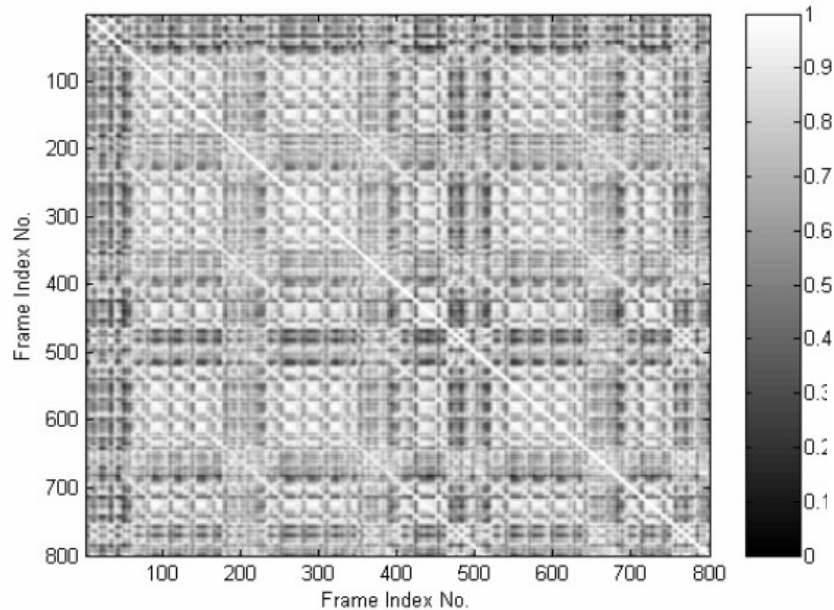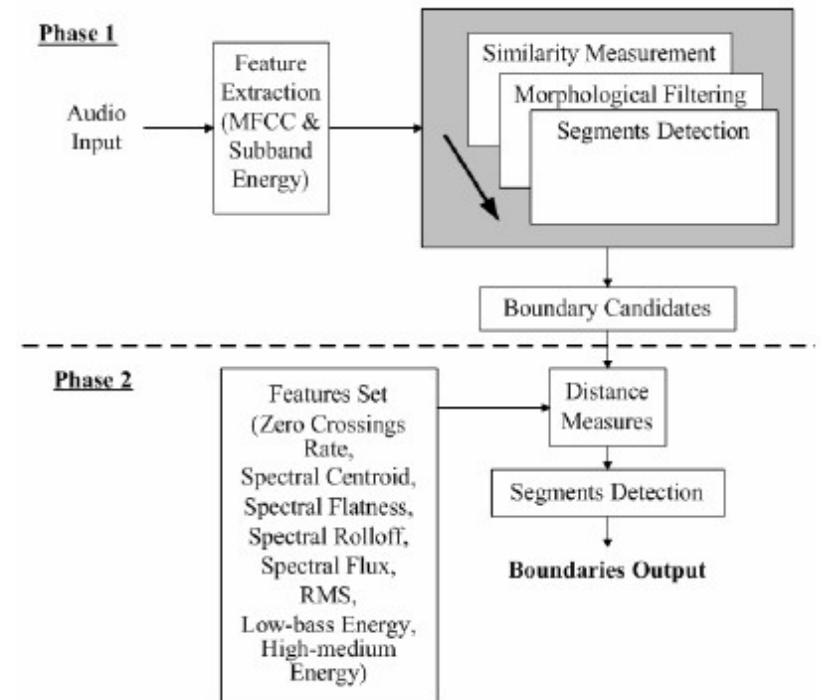


(Ong, 2006)

# Structure description



Figure 4.4. Two-dimensional similarity plot of The Beatles' song entitled *I'm a Loser*.
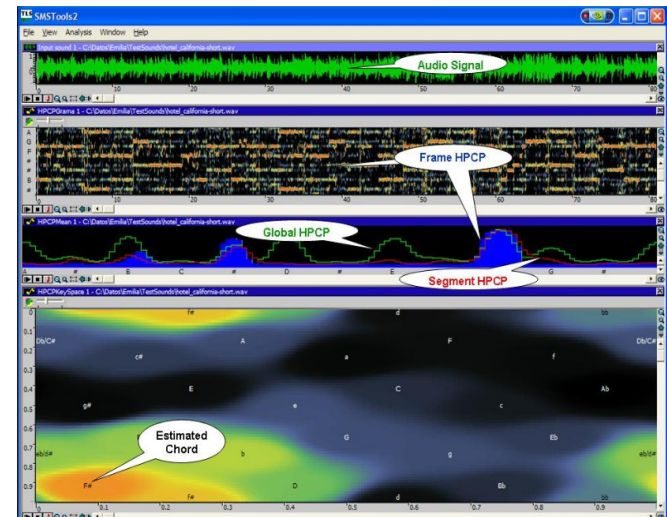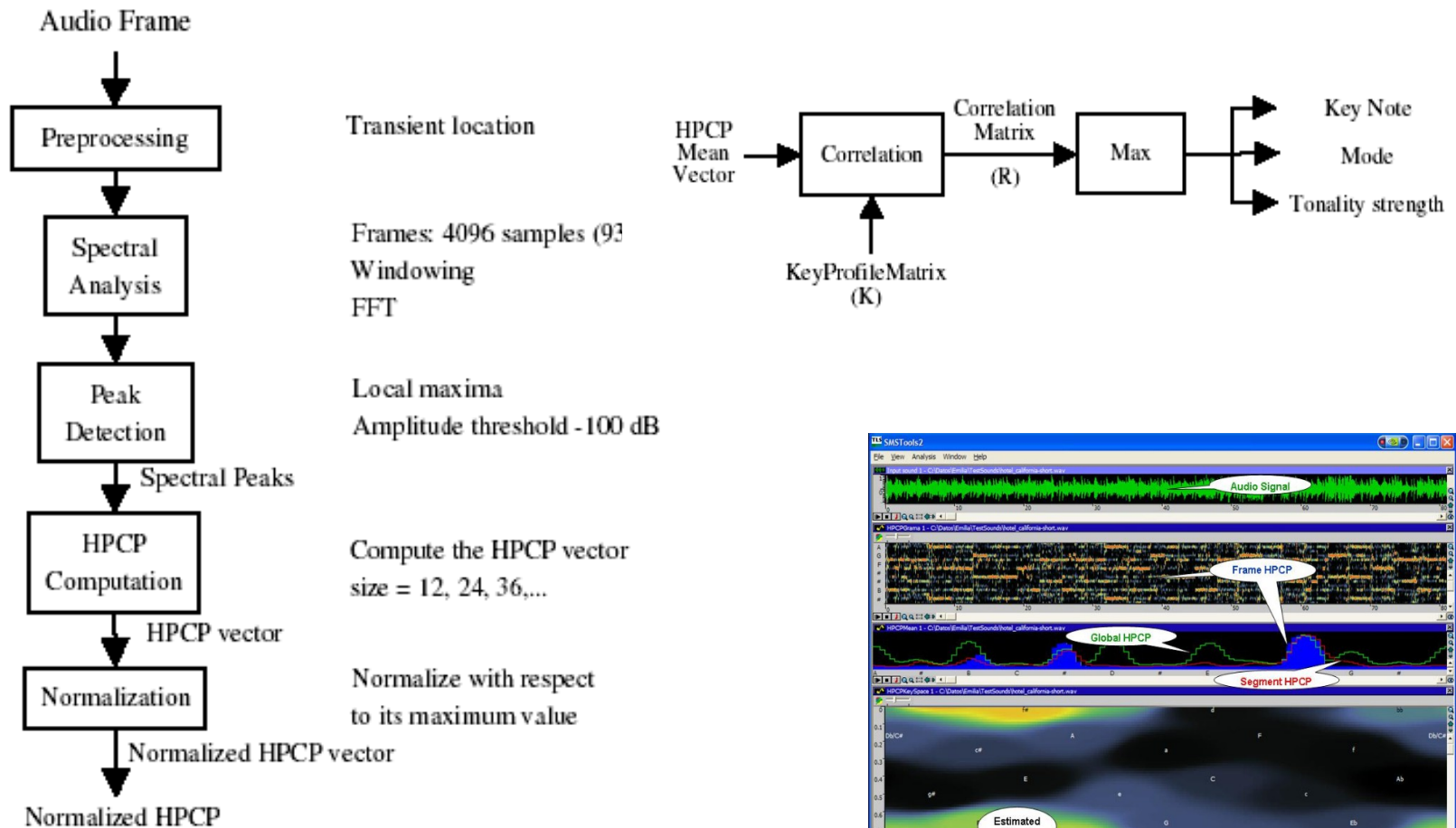
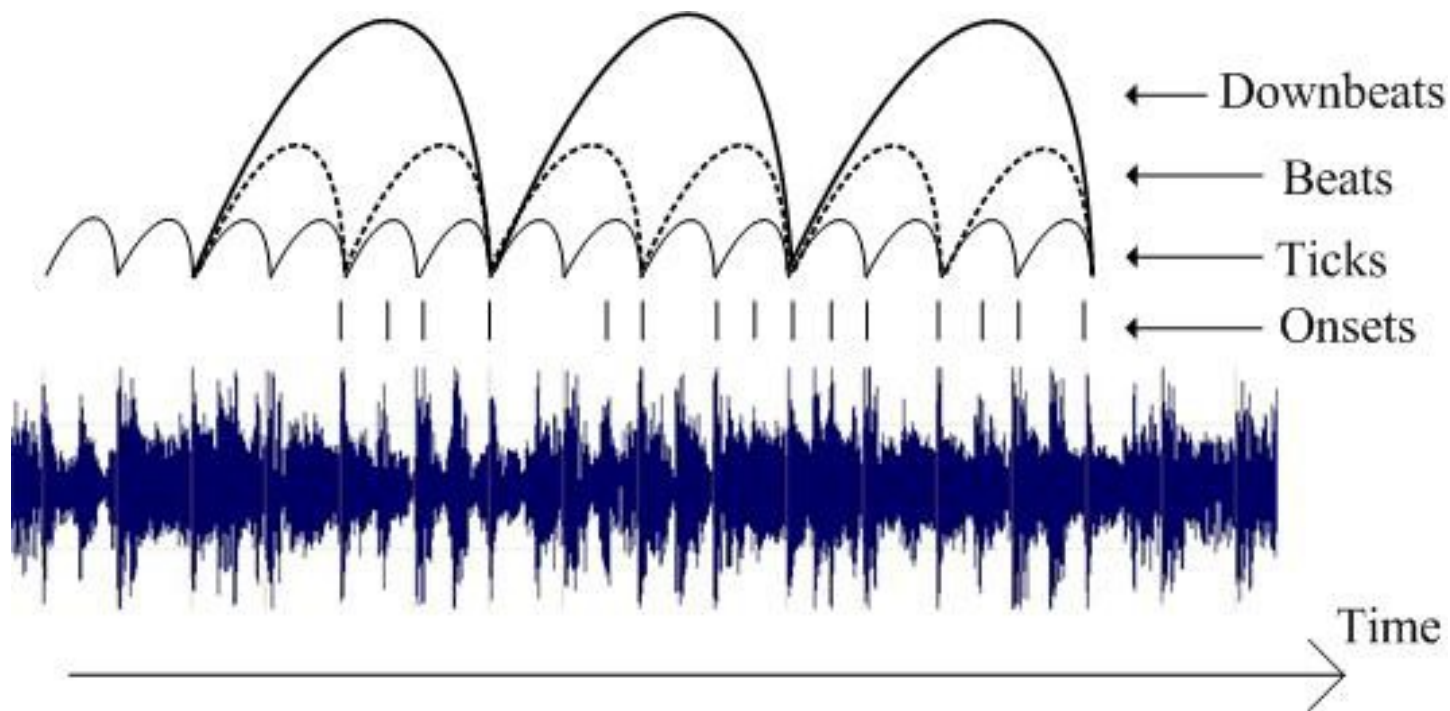(Ong, 2006)

# Tonal description

- Extract:
    - Melody (predominant melody or score)
    - Harmony (chords)
    - Key, modulations
- Much research is related to **automatic transcription** of music *(Klapuri PhD 2004)*
    - Fundamental frequency / Multipitch estimation (*de Cheveigné*)
    - Melody extraction (Predominant pitch, note segmentation)
    - Still unsolved, even for monophonic signals.
- Pitch class distribution of a piece
- Mid and high level features -> apply a tonal model / musical analysis (*Krumhansl, Leman, Temperley, ….*)
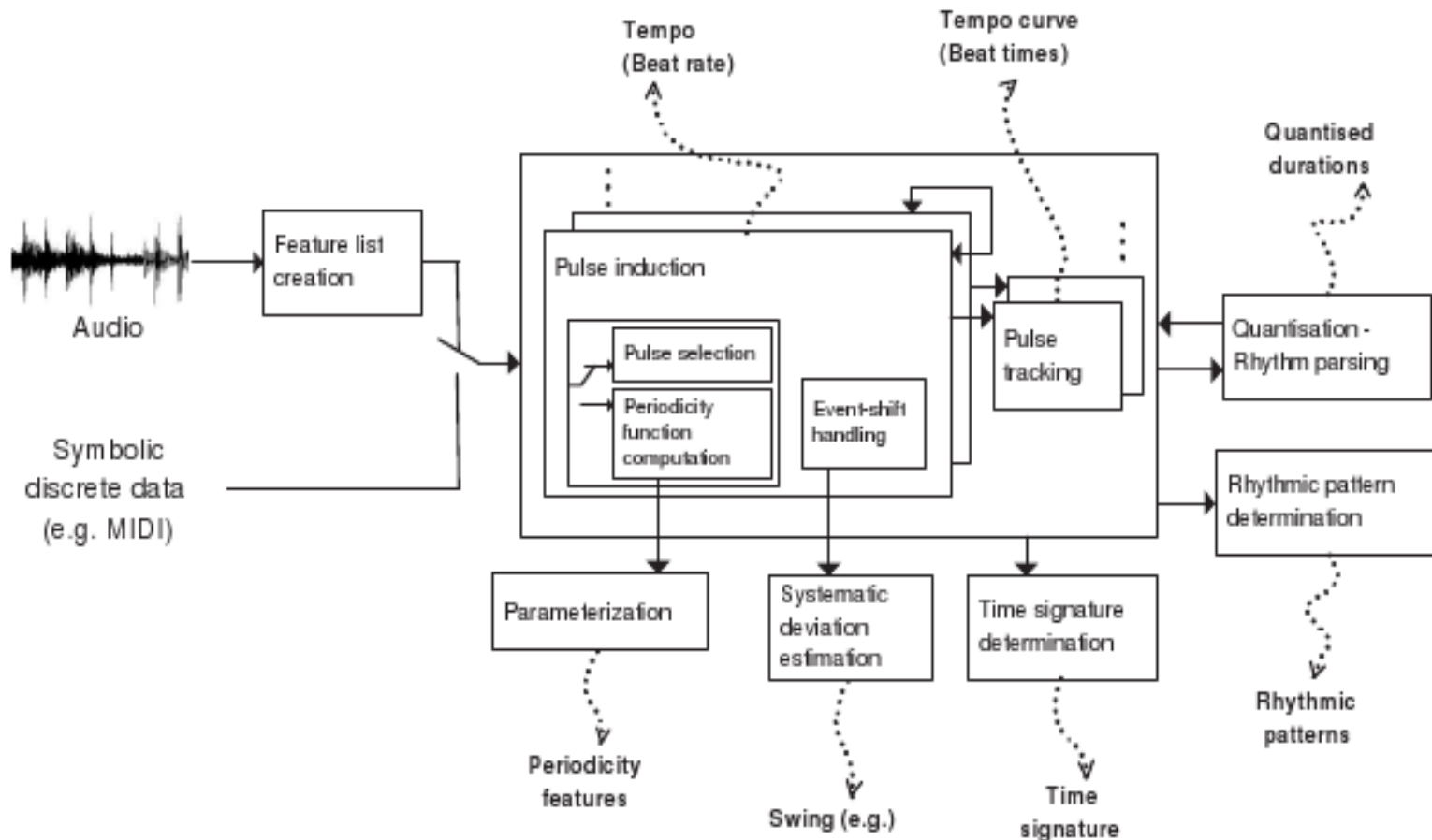
# Tonal description

# Rhythm description

Extraction of the metrical structure
of a piece



(Gouyon, 2005)

# Rhythm description



(Gouyon, 2005)

# References

- http://essentia.upf.edu

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

-

# Credits

All the slides of this presentation are released under an Attribution-Noncommercial-Share Alike license.